

# Optimisation:

De l'estimation paramétrique à l'apprentissage,  
une ballade entre théorie et pratique

S. Delprat

## Chapitre 2 – Optimisation sans contrainte

## Chapitre 2 – Optimisation sans contrainte

---

- 1) **Contexte**
- 2) Algorithme à base de Gradient
- 3) Méthode de Newton & quasi-Newton
- 4) Méthode « Régions de confiance » (Trust region)

## Contexte

---

Soit  $f$  une fonction convexe de classe  $\mathcal{C}^2$ .

On cherche le minimum (global) de la fonction  $f$ :  $\min_{x \in \mathbb{R}^n} f(x)$

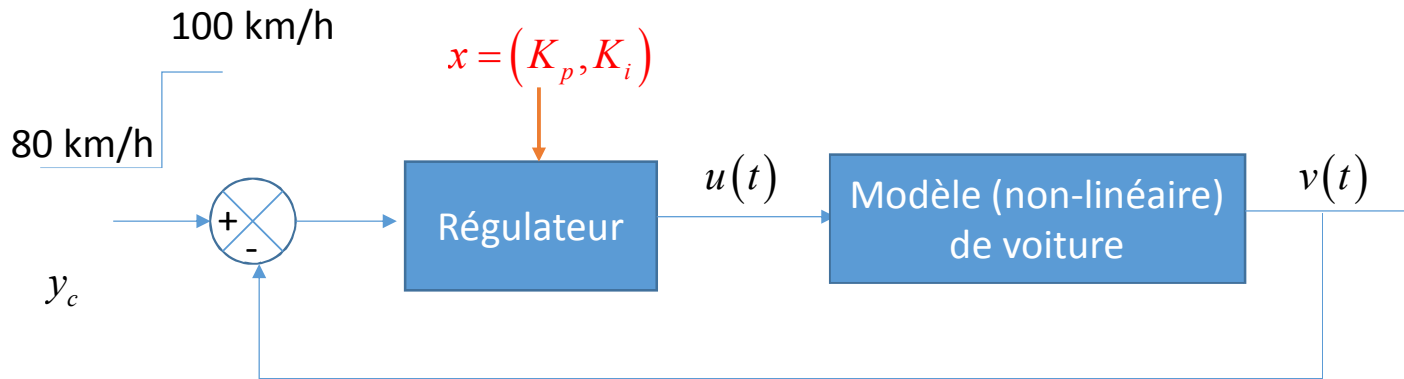
Nb: en pratique, si  $f$  n'est pas convexe, on trouvera un des minimum locaux

### Problème:

On peut évaluer la fonction  $f$  et son gradient mais on ne connaît pas la fonction  $f$ .

# Contexte

Problème: Trouver les paramètres d'un régulateur de vitesse.



$u(t)$  : couple délivré par le moteur thermique

$v(t)$  : vitesse du véhicule

Objectif: minimiser :

1. L'erreur de vitesse  $|y_c(t) - v(t)|$
2. La consommation d'énergie  $u(t)^2$

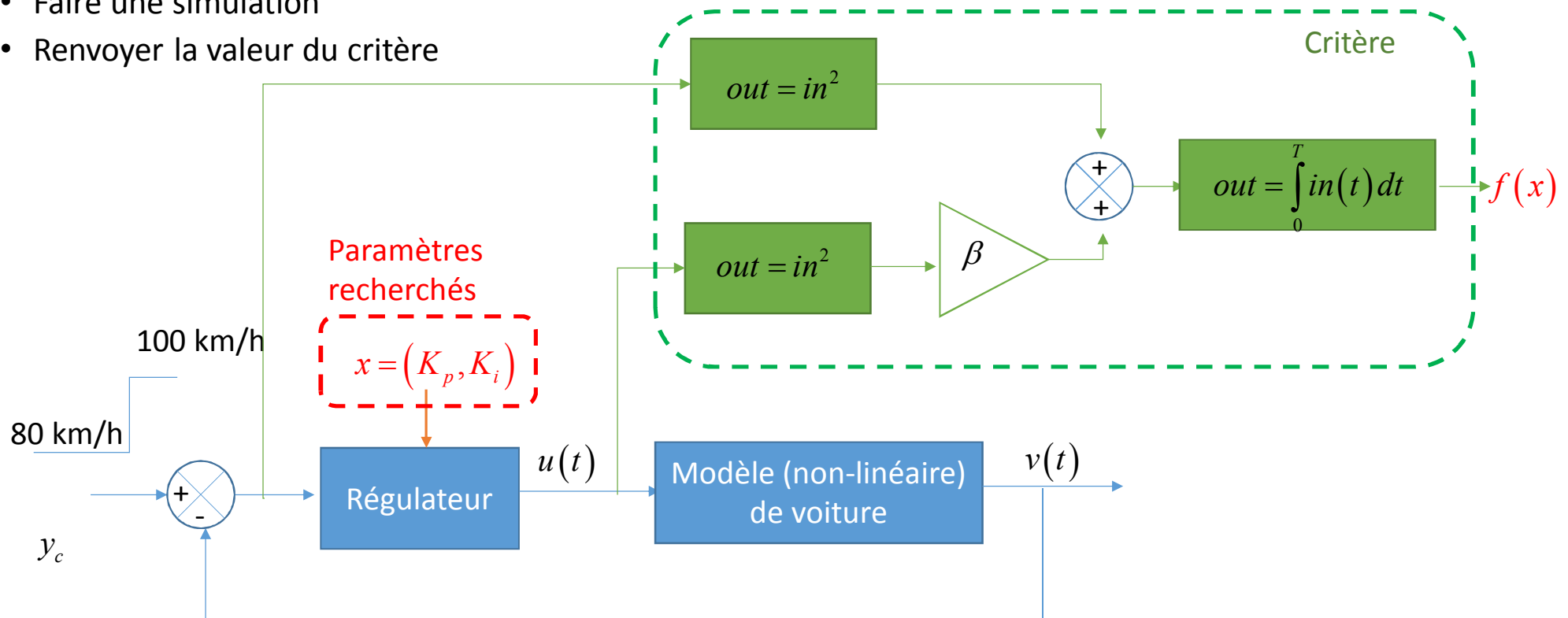
Compromis conso-perfo

Attention: on veut obtenir le minimum sur une durée donnée  $[0, T]$ : 
$$f(x) = \int_0^T \left[ (y_c(t) - v(t))^2 + \beta u(t)^2 \right] dt$$

# Contexte

**Evaluer la fonction  $f(x)$ , c'est:**

- Fixer les valeurs de  $x$
- Modifier les paramètres  $K_p, K_i$  du modèle en conséquence
- Faire une simulation
- Renvoyer la valeur du critère



## Contexte

---

Soit  $f$  une fonction convexe de classe  $\mathcal{C}^2$ .

On cherche le minimum (global) de la fonction  $f$ :  $\min_{x \in \mathbb{R}^n} f(x)$

Nb: en pratique, si  $f$  n'est pas convexe, on trouvera un des minimum locaux

### Problème:

On peut évaluer la fonction  $f$  et son gradient mais on ne connaît pas la fonction  $f$ .

### Solution:

On essaye de construire une suite de valeur qui converge vers le minimum global (ou local si  $f$  n'est pas convexe).

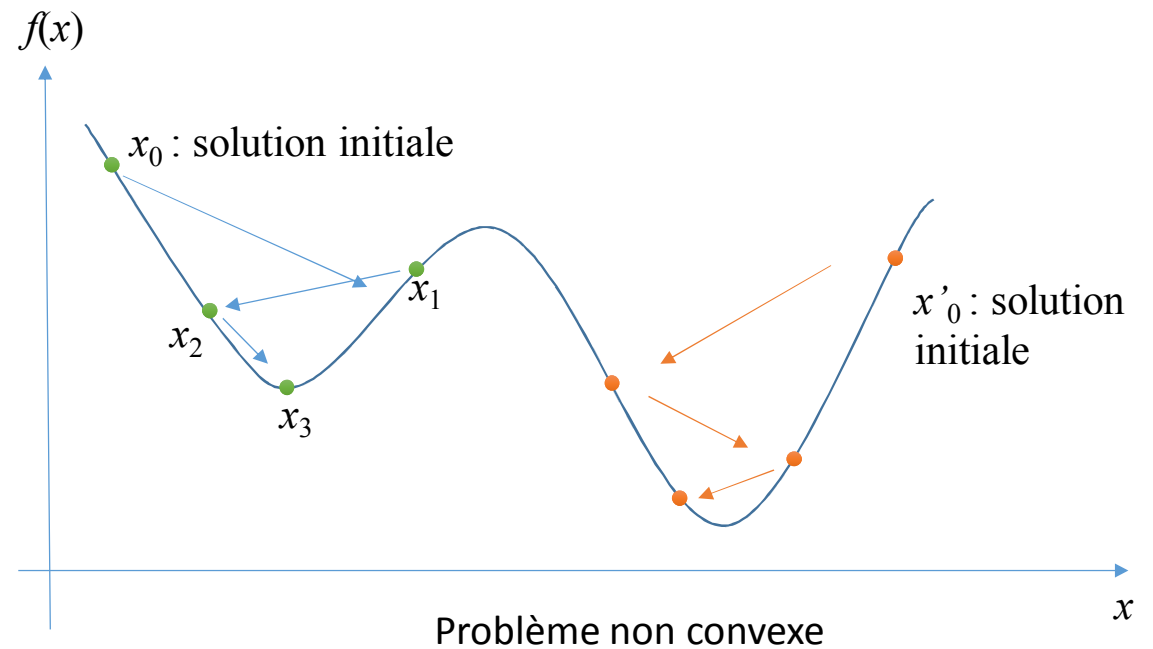
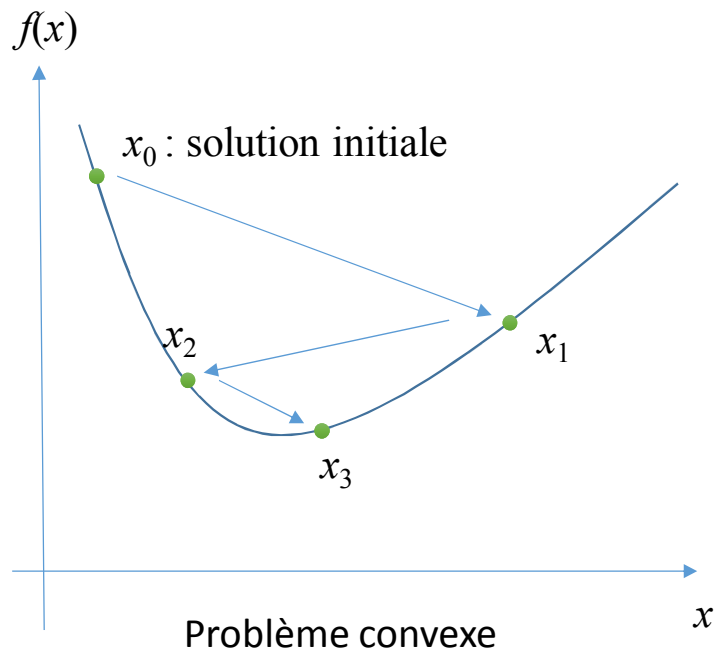
Problème : Il faut un point de départ.

# Contexte

Soit  $f$  une fonction convexe de classe  $\mathcal{C}^2$ .

On cherche le minimum (global) de la fonction  $f$ :  $\min_{x \in \mathbb{R}^n} f(x)$

Nb: en pratique, si  $f$  n'est pas convexe, on trouvera un des minimum locaux



# Préambule: Recherche d'un zéro d'une fonction

## • Méthode de bi-section

Soit une fonction  $g$  monotone (au moins localement)  $g : \mathbb{R} \rightarrow \mathbb{R}$

On cherche une racine  $g(x^*) = 0$

On connaît 2 bornes  $a^0$  et  $b^0$  encadrant la solution :  $g(a^0) < 0 < g(b^0)$

Algorithme

$k=0$

Répéter:

- On coupe l'intervalle en 2:  $c^k = \frac{a^k + b^k}{2}$

- Si  $g(c^k) > 0$  ALORS

$$a^{k+1} = c^k$$

$$b^{k+1} = b^k$$

SINON

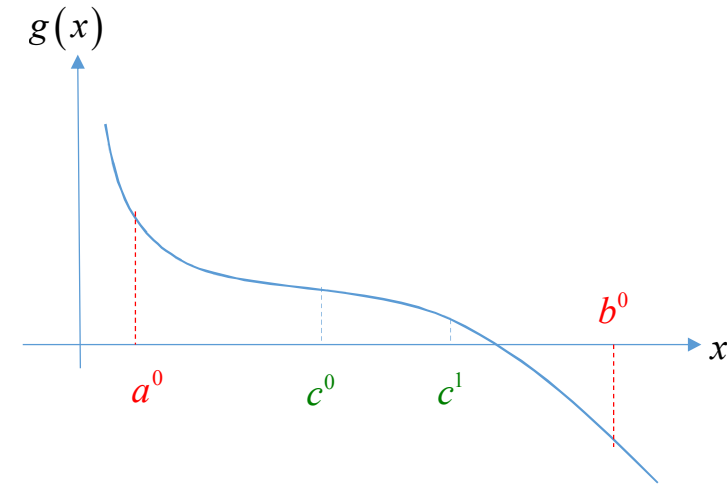
$$a^{k+1} = a^k$$

$$b^{k+1} = c^k$$

On garde l'intervalle qui contient la racine

Fin Si

Jusqu'à  $b^k - a^k < \varepsilon$





# Préambule: Recherche d'un zéro d'une fonction

## • Méthode de bi-section

```
clear all;
close all;
clc;

% Algorithme de bisection
f=@(x) (x-3).^3;

x=linspace(0,5,200);
figure;plot(x,f(x),'r')
grid on
hold on;
racine=bisection(f,-10,10,1e-5);
plot(racine,f(racine),'b*');

function c=bisection(f,a,b,eps)
ended=0;
while ended==0
    c=(a+b)/2;
    fc=f(c);
    if fc>0
        b=c;
    else
        a=c;
    end
    ended=(b-a)<eps;
end
end
```

### Savoir-faire n°1 : Calcul Matlab

Savoir trouver avec une précision donnée le zéro d'une fonction monovariante sur un intervalle donné

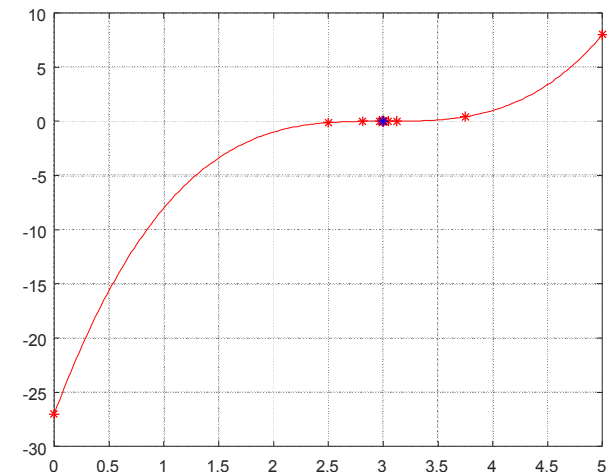
Intérêt :

- Simple à mettre en œuvre
- Pas besoin de connaître la dérivée de  $f$

Inconvénient:

- Convergence lente

$$f(x) = (x-3)^3$$



## Chapitre 2 – Optimisation sans contrainte

---

- 1) Contexte
- 2) **Algorithme à base de Gradient**
- 3) Méthode de Newton & quasi-Newton
- 4) Méthode « Régions de confiance » (Trust region)

# Direction de recherche

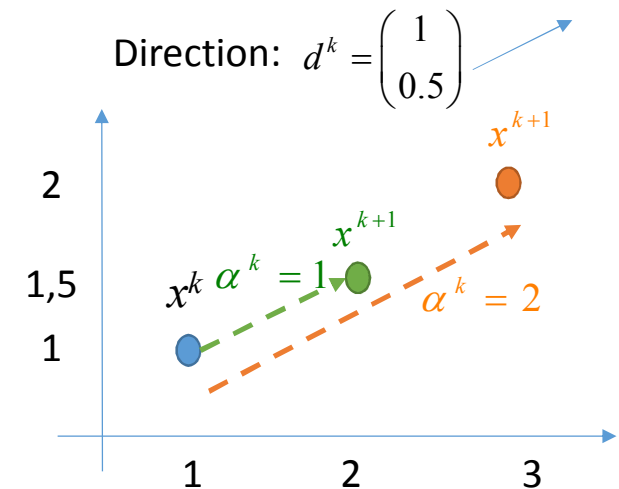
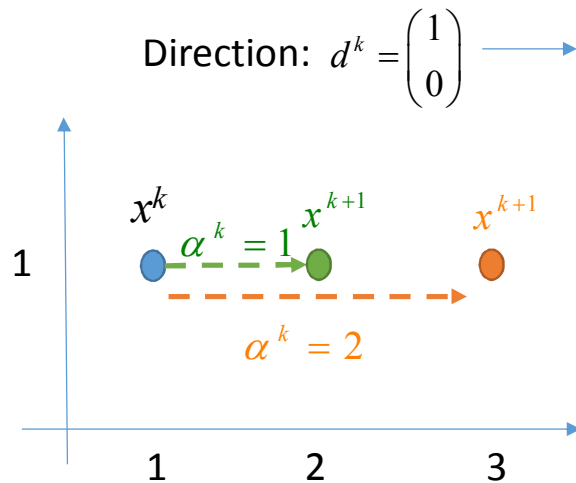
Construction d'une série qui converge vers un minimum local de  $f$ :

$x_0$  : point de départ, fournit par l'utilisateur ou aléatoire

$$x^{k+1} = x^k + \alpha^k d^k$$

$d^k \in \mathbb{R}^{n,1}$  : Direction de recherche

$\alpha^k \in \mathbb{R}$  : Taille du pas de recherche



## Vocabulaire: direction de descente

Soit  $x \in \mathbb{R}^n$ ,  $w \in \mathbb{R}^n \setminus \{0\}$  est une direction de descente en  $x$  s'il existe  $\alpha_0 > 0$  tel que

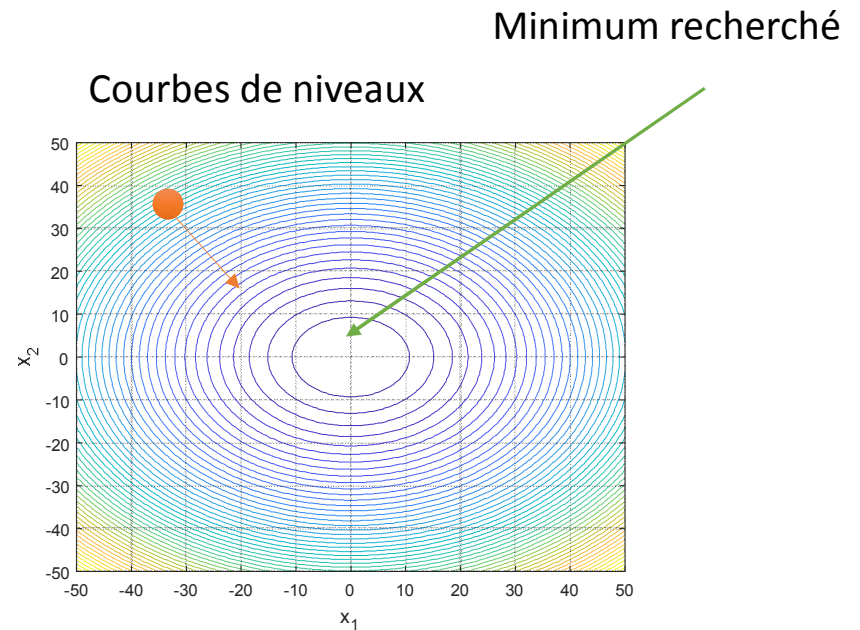
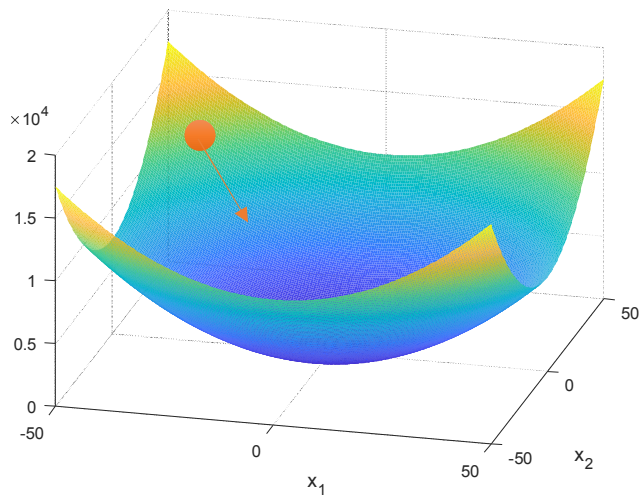
$$f(x + \alpha d) \leq f(x), \quad \forall \alpha \in [0, \alpha_0]$$

Soit  $x \in \mathbb{R}^n$ ,  $w \in \mathbb{R}^n \setminus \{0\}$  est une direction de descente *stricte* en  $x$  s'il existe  $\alpha_0 > 0$  tel que

$$f(x + \alpha d) < f(x), \quad \forall \alpha \in ]0, \alpha_0[$$

$\alpha \in \mathbb{R}^{+*}$   
 $d \in \mathbb{R}^{n,1}$

Donc :  $d$  donne la direction et  $\alpha$  la « taille » du pas

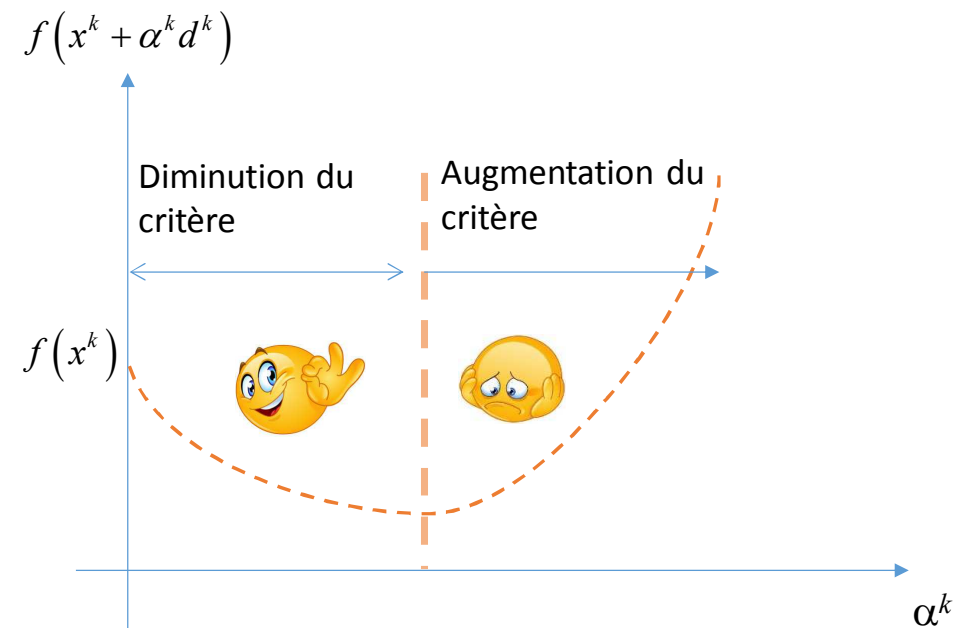
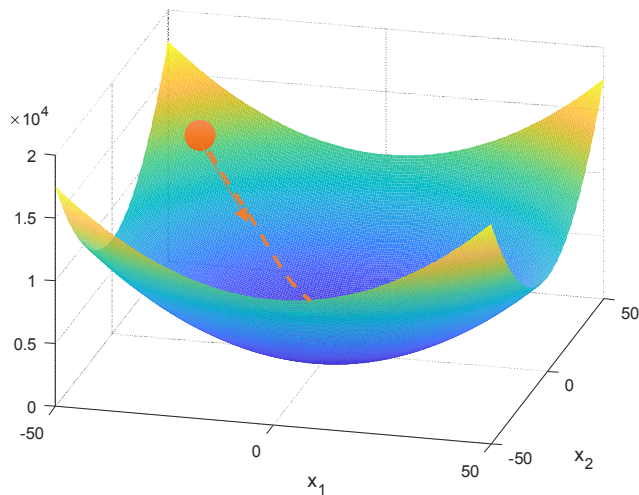


# Vocabulaire: direction de descente

Pour une direction de descente  $d^k$ :

- Certaines valeurs de pas  $\alpha^k$  trop grandes peuvent conduire à  $x^{k+1} > x^k$
- Pour une fonction convexe, il existe toujours une valeur de pas  $\alpha^k$  qui garantit  $x^{k+1} < x^k$  (tant que  $x^k$  n'est pas un minimum local)

Donc :  $d$  donne la direction et  $\alpha$  la « taille » du pas



# Algorithme générique

Paramètres de l'algorithme:

Solution initiale:  $x^0$

$k = 0$

Tant que *Critère d'arrêt* non satisfait

Calculer une direction de descente  $d^k$

Calculer un pas  $\alpha^k$  pour la direction de descente  $d^k$

Mettre à jour la solution :  $x^{k+1} = x^k + \alpha^k d^k$

Passage à l'itération suivante:  $k = k + 1$

Ce qui différencie les algorithmes:

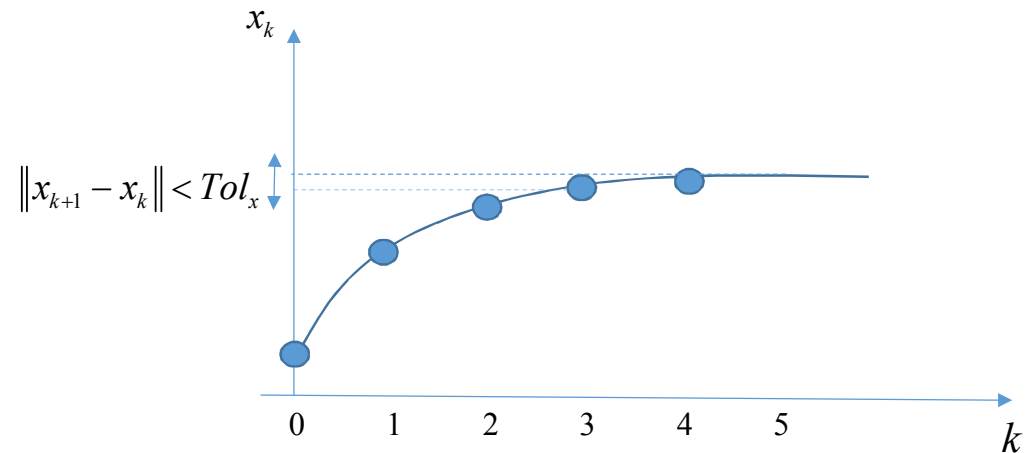
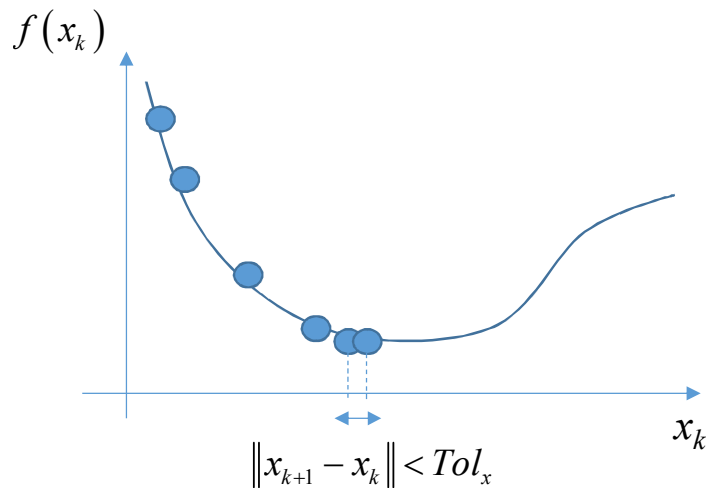
- Le choix de la direction de descente  $d^k$
- Le choix du pas  $\alpha^k$

# Algorithme générique : Critère d'arrêt

## Critères d'arrêts:

En pratique les solveurs permettent de spécifier plusieurs critères. L'algorithme s'arrête dès que l'un d'entre eux est satisfait

- Nombre d'itérations maximal est atteint
- Nombre d'évaluation maximal de la fonction *f atteint*.
- La solution est trouvée avec la bonne précision:  $\|x_{k+1} - x_k\| < Tol_x$   
=> Permet de maîtriser le nombre de décimales significatives

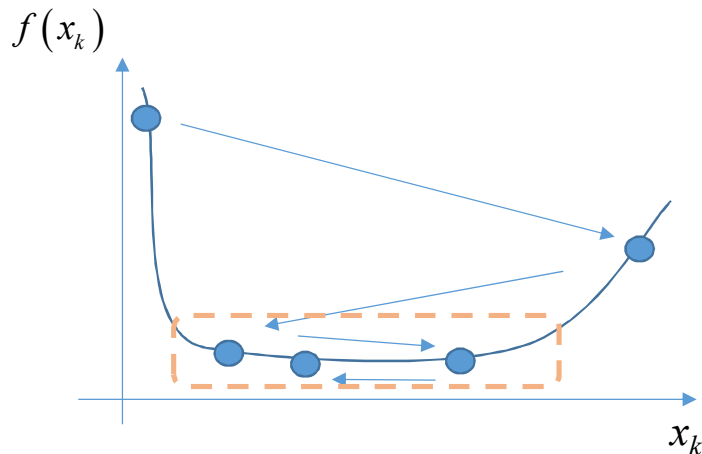


# Algorithme générique : Critère d'arrêt

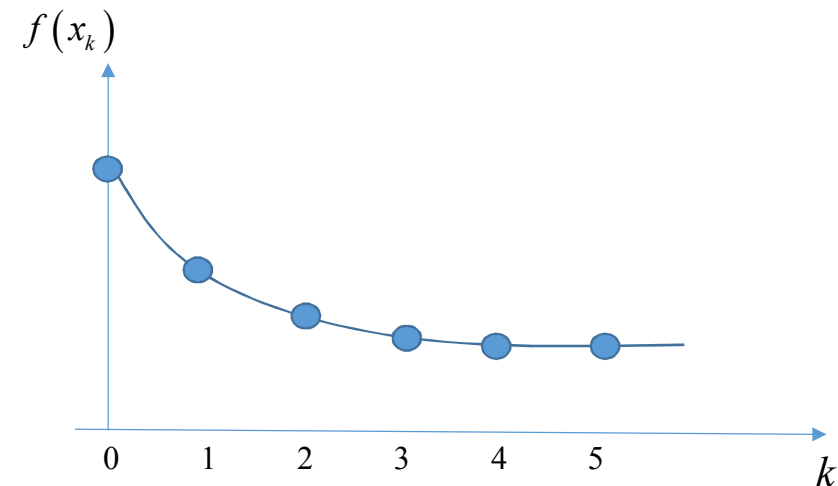
## Critères d'arrêts:

En pratique les solveurs permettent de spécifier plusieurs critères. L'algorithme s'arrête dès que l'un d'entre eux est satisfait

- L'algorithme n'améliore plus la valeur de la fonction:  $\|f(x_{k+1}) - f(x_k)\| < Tol_f$   
=> Permet d'éviter des itérations inutiles si la fonction  $f$  est très « plate » autour du minimum local



Toutes ces valeurs de  $x_k$  correspondent quasiment la même valeur du critère





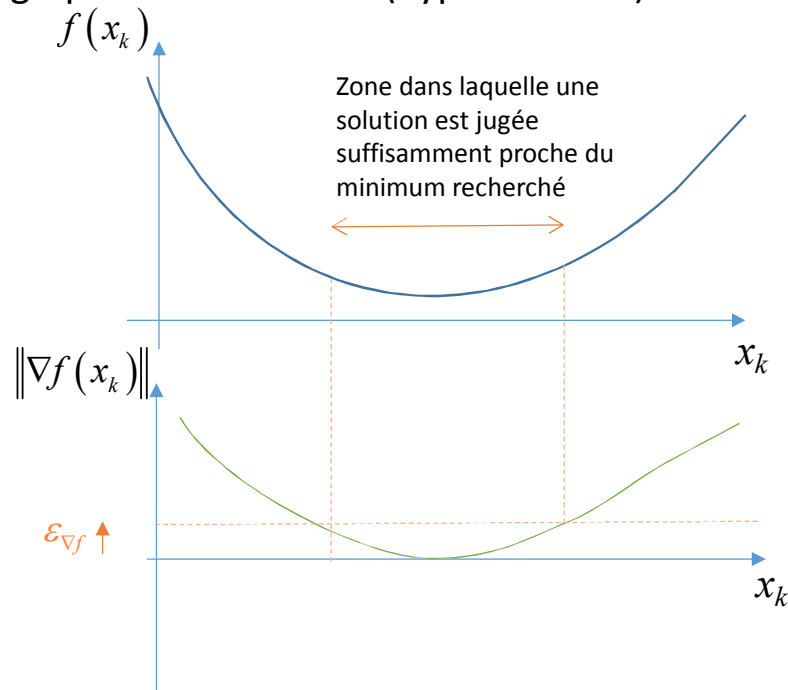
# Algorithme générique : Critère d'arrêt

## Critères d'arrêts:

En pratique les solveurs permettent de spécifier plusieurs critères. L'algorithme s'arrête dès que l'un d'entre eux est satisfait

- La condition d'optimalité au premier ordre est quasiment vérifiée:  $\|\nabla f(x_k)\| < \varepsilon_{\nabla f}$

=> Le graphe de la fonction (hyper-surface) est suffisamment « plate »



# Algorithme générique : Direction de descente

Paramètres de l'algorithme:

Solution initiale:  $x^0$

$k = 0$

Tant que *Critère d'arrêt* non satisfait

Calculer une direction de descente  $d^k$

Calculer un pas  $\alpha^k$  pour la direction de descente  $d^k$

Mettre à jour la solution :  $x^{k+1} = x^k + \alpha^k d^k$

Passage à l'itération suivante:  $k = k + 1$

## Algorithme générique : Direction de descente

On cherche une direction de descente  $d^k$  avec  $x^{k+1} = x^k + \alpha^k d^k$

On utilise le développement de Taylor en  $x^k$  pour une valeur suffisamment petite du pas  $\alpha^k$

$$f(x^{k+1}) = f(x^k + \alpha^k \textcolor{red}{d}^k) = f(x^k) + \nabla f(x^k)^T \alpha^k \textcolor{red}{d}^k + O(\alpha^k d^k)$$

$$f(x^{k+1}) - f(x^k) = \nabla f(x^k)^T \alpha^k \textcolor{red}{d}^k + O(\alpha^k d^k)$$

Pour avoir une direction de descente  $f(x^{k+1}) - f(x^k) < 0$

$$\begin{array}{l} \nabla f(x^k)^T \alpha^k d^k + O(\alpha^k d^k) < 0 \\ \nabla f(x^k)^T d^k + \frac{O(\alpha^k d^k)}{\alpha^k} < 0 \end{array} \quad \begin{array}{l} \downarrow \\ \alpha^k > 0 \\ \leftarrow \end{array} \quad \lim_{\alpha^k \rightarrow 0^+} \frac{O(\alpha^k d^k)}{\alpha^k} = 0$$

Finalement, on a une condition pour obtenir une direction de descente :

$$\nabla f(x^k)^T \textcolor{red}{d}^k < 0$$

## Direction de descente : Méthode de la plus forte pente

On cherche une direction de descente:  $x^{k+1} = x^k + \alpha^k d^k$

Intuition : descendre dans le sens de la plus forte pente:  $d^k = -\nabla f(x^k)$

$$x^{k+1} = x^k - \alpha^k \nabla f(x^k)$$

C'est évidemment une direction de descente:

$$\nabla f(x^k)^T d^k = -\nabla f(x^k)^T \nabla f(x^k) < 0$$

On montre que c'est la plus forte direction de descente.

Une alternative serait de choisir  $d^k = -M \nabla f(x^k)$  Avec  $M$  définie positive

On parle alors de méthode à base de gradient

## Direction de descente : Méthode de la plus forte pente

On cherche une direction de descente:  $x^{k+1} = x^k + \alpha^k d^k$

Intuition : descendre dans le sens de la plus forte pente:  $d^k = -\nabla f(x^k)$

Problème : la norme du gradient est variable.

Solution : normalisation de la direction (facilitera la recherche de  $\alpha^k$ )

$$d^k = -\frac{\nabla f(x^k)}{\|\nabla f(x^k)\|}$$

On obtient toujours la même direction mais la norme du vecteur est toujours unitaire.

# Algorithme générique : calcul du pas $\alpha^k$

Paramètres de l'algorithme:

Solution initiale:  $x^0$

$k = 0$

Tant que *Critère d'arrêt* non satisfait

Calculer une direction de descente  $d^k$

Calculer un pas  $\alpha^k$  pour la direction de descente  $d^k$

Mettre à jour la solution :  $x^{k+1} = x^k + \alpha^k d^k$

Passage à l'itération suivante:  $k = k + 1$

# Exemple académique : fonction quadratique

$$f(x,y) = (x \ y) \begin{pmatrix} 1 & 0.5 \\ 0.5 & 2 \end{pmatrix} \begin{pmatrix} x \\ y \end{pmatrix}$$

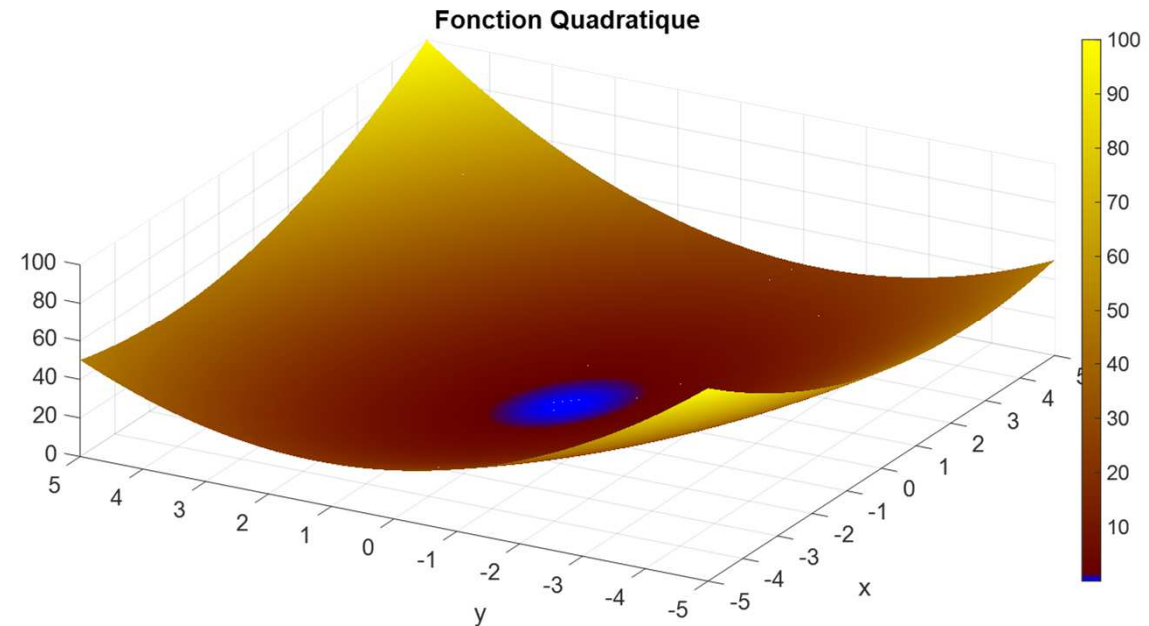
```
clear all;
close all;
clc;

H=[1 0.5;
   0.5 2];
f=@(x) [x(1) x(2)]*H*[x(1);x(2)];
xmin=-5;
xmax=5;
ymin=-5;
ymax=5;
Fname='Fonction Quadratique';

% Définition de la colormap
Delta1=(0:10)'/10;
Delta2=(0:750)'/750;
map=[0.4*Delta1 0*Delta1 1-Delta1;
     0.4+0.6*Delta2 1*Delta2 0*Delta2];

% Calcul de la fonction
X=linspace(xmin,xmax,61);
Y=linspace(ymin,ymax,62);
Z=zeros(length(Y),length(X));
for i=1:length(X)
    for j=1:length(Y)
        Z(j,i)=f([X(i);Y(j)]);
    end
end

% Affichage
surf(X,Y,Z,'EdgeColor','none');
xlabel('x');
ylabel('y');
title(Fname)
colorbar;
colormap(map)
hold on;
```



# Exemple académique : fonction quadratique

On utilise un pas  $\alpha$  constant:

```
clear all;
close all;
clc;
H=[1 0.5;
   0.5 2];

f=@(x) [x(1) x(2)]*H*[x(1);x(2)];
df=@(x) H*[x(1);x(2)]+H'*[x(1);x(2)];
x0=[5;-5];
xmin=-5;
xmax=5;
ymin=-5;
ymax=5;
Fname='Quadratique';

% Définition de la colormap
Delta1=(0:10)'/10;
Delta2=(0:750)'/750;
map=[0.4*Delta1 0*Delta1 1-Delta1;
     0.4+0.6*Delta2 1*Delta2 0*Delta2];

% Criteres d'arret
Tolf=1e-6;
Tolx=1e-9;
IterMax=100;
```

Gradient calculé “manuellement”

```
% Affichage du contour
figure;
X=linspace(xmin,xmax,61);
Y=linspace(ymin,ymax,62);
Z=zeros(length(Y),length(X));
for i=1:length(X)
    for j=1:length(Y)
        Z(j,i)=f([X(i);Y(j)]);
    end
end
surf(X,Y,Z,'EdgeColor','none');
xlabel('x');
ylabel('y');
title(Fname)
colorbar;
colormap(map)
hold on;
```

```
% Algorithme
ended=0;
k=1;
fxkp1=f(x0);
x=x0;
```

Pas fixe

```
while ended==0
    dk=-df(x(:,k));
    dk=dk/norm(dk);
    alpha=0.3;
    x(:,k+1)=x(:,k)+alpha*dk;
    fxk=fxkp1;
    fxkp1=f(x(:,k+1));

    plot3(x(1,k+1),x(2,k+1),fxkp1,'m');
    plot3([x(1,k+1) x(1,k)], [x(2,k+1) x(2,k)], [fxkp1 fxk], 'g');

    ended=(k>IterMax) || ...
        norm(fxkp1-fxk)<Tolf || ...
        norm(x(:,k)-x(:,k+1))<Tolx;
    k=k+1;
end

fprintf('Valeur finale : x=[%.2f,%.2f] \n',x(1,k),x(2,k))
if k>IterMax
    fprintf(' Iteration maximale atteinte\n');
end
if norm(fxkp1-fxk)<Tolf
    fprintf('Plus d\'amélioration de f\n');
end
if norm(x(:,k)-x(:,k+1))<Tolx
    fprintf('plus d\'amélioration de x\n');
end
```

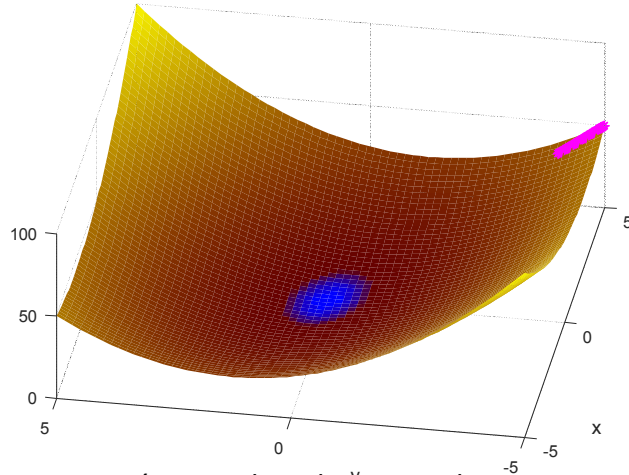


# Exemple académique : fonction quadratique

On utilise un pas alpha constant:

$$\alpha = 0.01$$

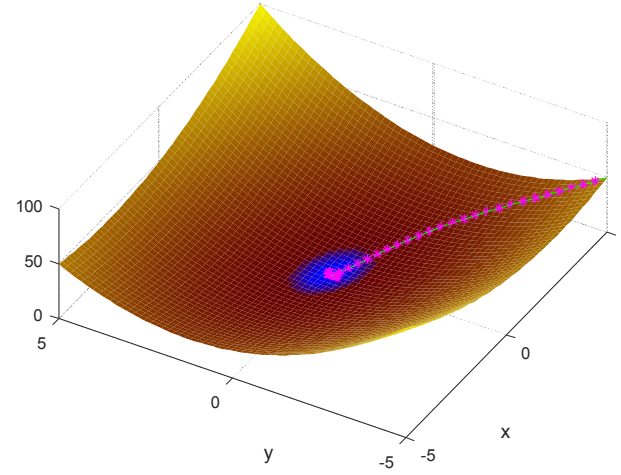
Fonction quadratique



102 itérations boucle principale  
102 appels de f

$$\alpha = 0.3$$

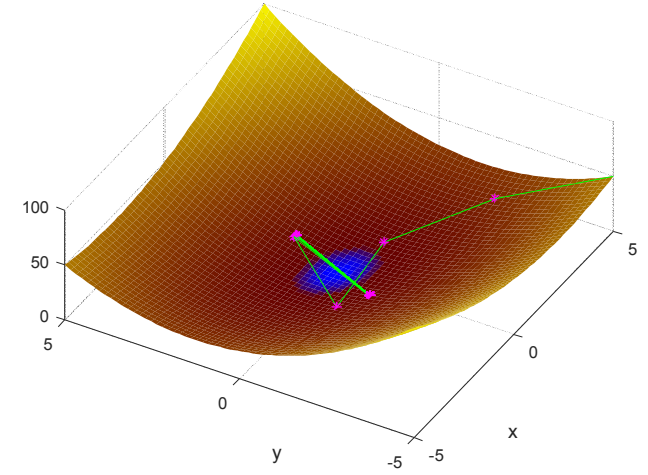
Fonction quadratique



102 itérations boucle principale  
102 appels de f

$$\alpha = 3$$

Fonction quadratique



102 itérations boucle principale  
102 appels de f

Problème: la valeur du « pas »  $\alpha$  est difficile à choisir, même pour une fonction aussi simple qu'une quadratique.  
=> On peut montrer que pour beaucoup d'exemples l'algorithme ne convergera pas ou convergera lentement.

Intuitivement, il faut adapter le pas à la géométrie locale de la fonction.

## Algorithme générique : calcul du pas $\alpha^k$

Le problème à résoudre est le choix du pas  $\alpha^k$  qui minimise le cout:

$$\alpha^k = \arg \min_{\alpha > 0} f(x^k + \alpha d^k)$$

Comme  $\alpha^k \in \mathbb{R}^{+*}$ , le problème est potentiellement « simple » à résoudre (car scalaire).

Plusieurs approches envisageables:

- Résolution exacte: possible que pour des cas particuliers.

En pratique souvent impossible à résoudre pour un problème réaliste.  
i.e. convient bien si  $f$  est une forme quadratique

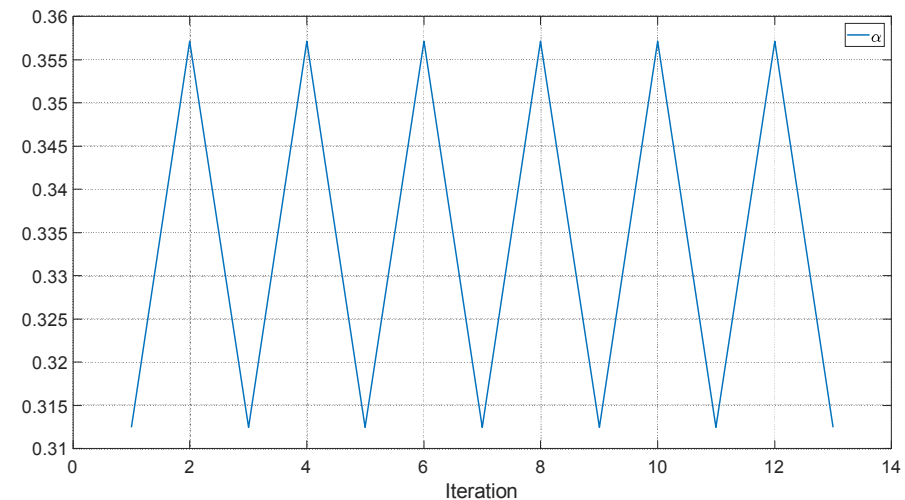
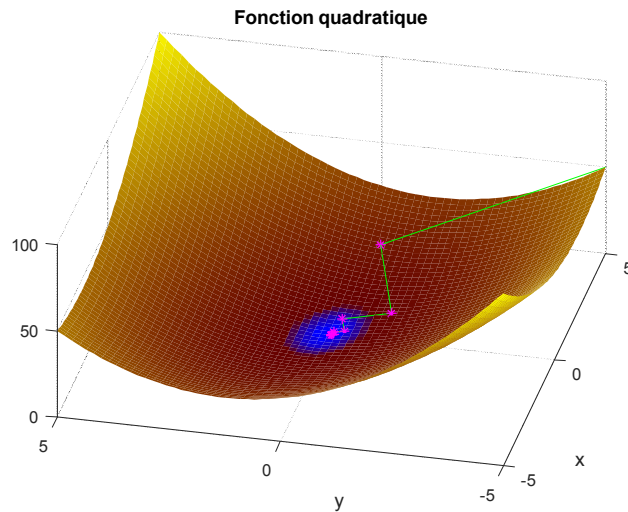
=> Cas d'une fonction quadratique  $f(x) = x^T Hx + G^T x + A$

Pour: 
$$x^{k+1} = x^k - \alpha^k \frac{\nabla f(x^k)}{\|\nabla f(x^k)\|}$$

Le pas optimal est: 
$$\alpha^k = 0.5 \frac{(d^k)^T d^k}{(d^k)^T H d^k} \|\nabla f(x^k)\|$$

## Exemple académique : fonction quadratique

On utilise un pas  $\alpha^k = 0.5 \frac{(d^k)^T d^k}{(d^k)^T H d^k} \|\nabla f(x^k)\|$  optimal solution de :  $\alpha^k = \arg \min_{\alpha > 0} f(x^k + \alpha d^k)$



Les directions de recherche successives sont alors orthogonales

Problème : en pratique le pas optimal n'est pas calculable.

=> Idée : s'assurer que la valeur de la fonction diminue entre 2 itérations

## Algorithme générique : calcul du pas $\alpha^k$ — *Algorithme Backtracking*

Le problème à résoudre est le choix du pas  $\alpha^k$  qui minimise le cout:  $\alpha^k = \arg \min_{\alpha > 0} f(x^k + \alpha d^k)$

- Algorithme de « Backtracking »

Comme  $d^k$  est une direction de descente, il suffit de prendre un pas  $\alpha$  assez petit pour avoir

$$f(x^k + \alpha d^k) < f(x^k)$$

Idée : partir d'un pas initial et le diminuer jusqu'à ce que la condition soit satisfaite.

### Algorithme au pas $k$ :

Paramètres :  $\alpha_0^k > 0, \beta \in ]0,1[$

$i = 0$

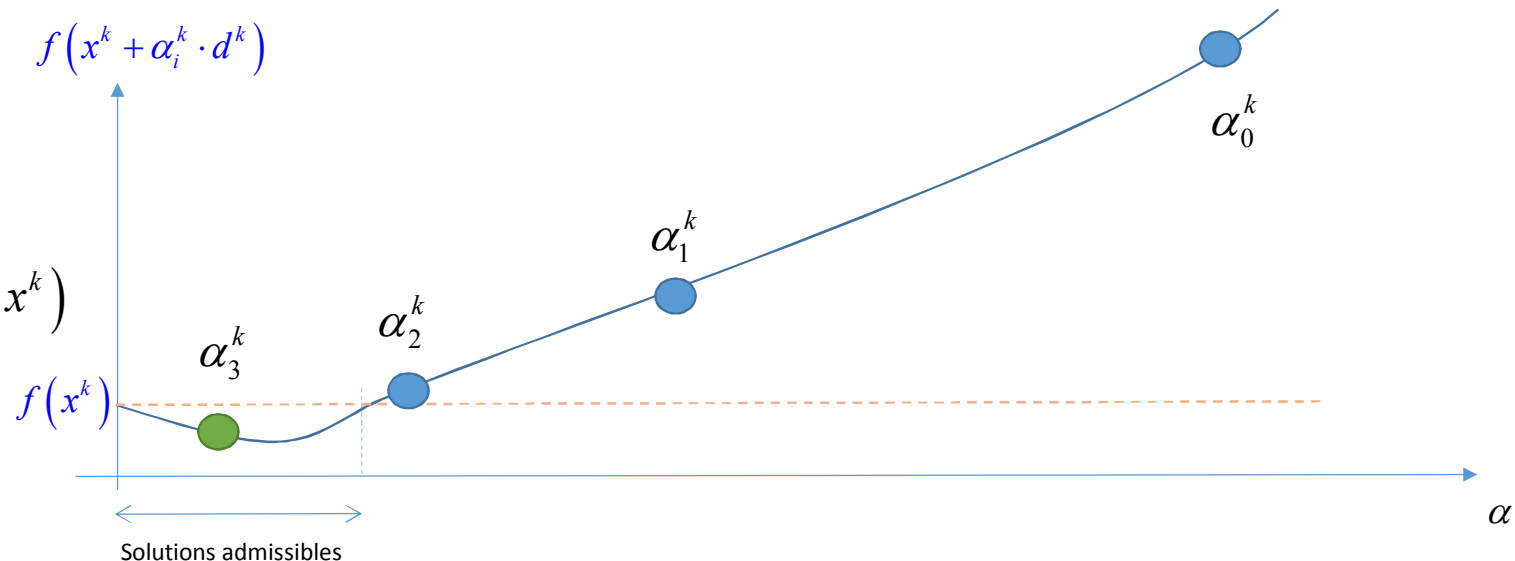
Tant que  $f(x^k + \alpha_i^k d^k) \geq f(x^k)$

$$\alpha_{i+1}^k = \alpha_i^k \cdot \beta$$

$i = i + 1$

Fin Tant que

$$\alpha^k = \alpha_i^k$$



# Algorithme générique : calcul du pas $\alpha^k$ — *Algorithme Backtracking*

Modifications du programme précédent:

```
% Paramètres backtracking
alpha0=1;
beta=0.9;

while ended==0
    dk=-df(x(:,k));

    alpha(k)=BackTracking(f,x(:,k),dk,alpha0,beta);

    x(:,k+1)=x(:,k)+alpha(k)*dk;
    fxk=fxkp1;
    fxkp1=f(x(:,k+1));

    plot3(x(1,k+1),x(2,k+1),fxkp1,'m*');
    plot3([x(1,k+1) x(1,k)], [x(2,k+1) x(2,k)], [fxkp1 fxk], 'g');

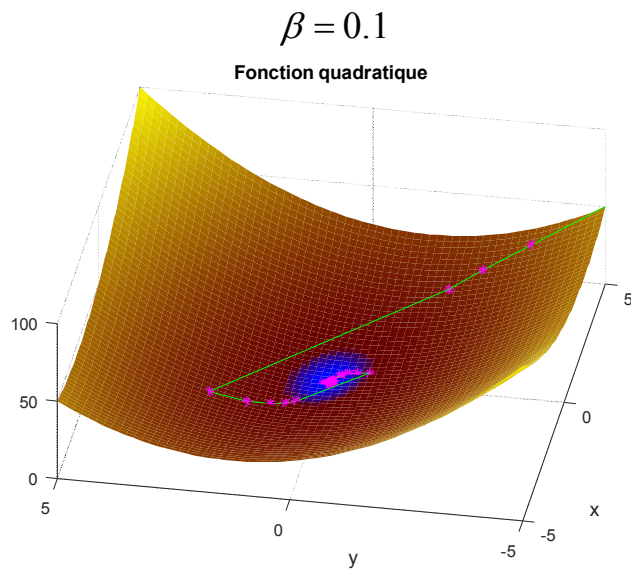
    ended=(k>IterMax) || ...
        norm(fxkp1-fxk)<Tol1f || ...
        norm(x(:,k)-x(:,k+1))<Tolx;
    k=k+1;
end

function alpha=BackTracking(ff,xk,dk,alpha0,beta)
fxk=f(xk);
alpha=alpha0;
while fxk<f(xk+alpha*dk)
    alpha=alpha*beta;
end
end
```

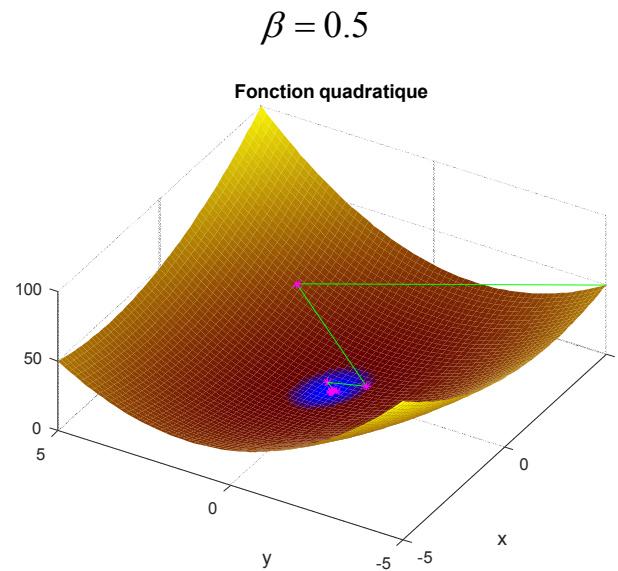
NB: il est possible d'optimiser le programme pour économiser des appels à la fonction f...

## Exemple académique : fonction quadratique

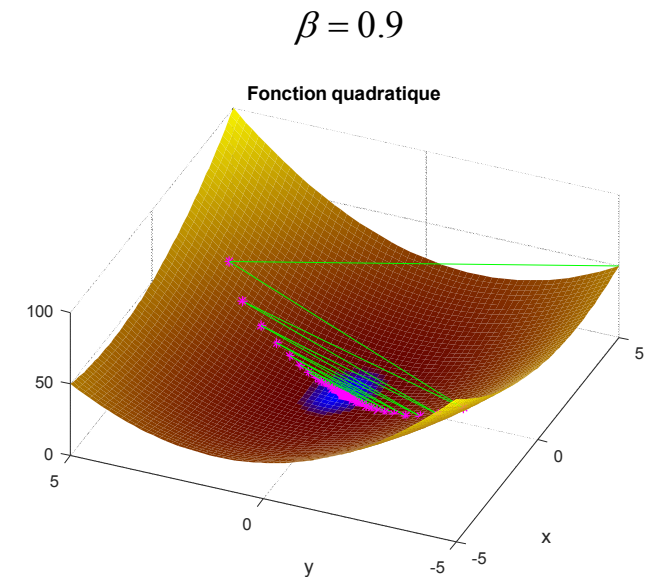
On utilise un pas calculé avec l'algorithme de Backtracking avec:  $\alpha_0^k = 1$



37 itérations boucle principale  
137 appels de f



12 itérations boucle principale  
49 appels de f



79 itérations boucle principale  
855 appels de f

Problème: l'algorithme de backtracking s'arrête à la première amélioration de f  
=> Peut être que des pas plus petits seraient plus efficace.

## Algorithme générique : calcul du pas $\alpha^k$ – *Algorithme Backtracking*

Le problème à résoudre est le choix du pas  $\alpha^k$  qui minimise le cout:  $\alpha^k = \arg \min_{\alpha > 0} f(x^k + \alpha d^k)$

- **Algorithme de « Backtracking »**

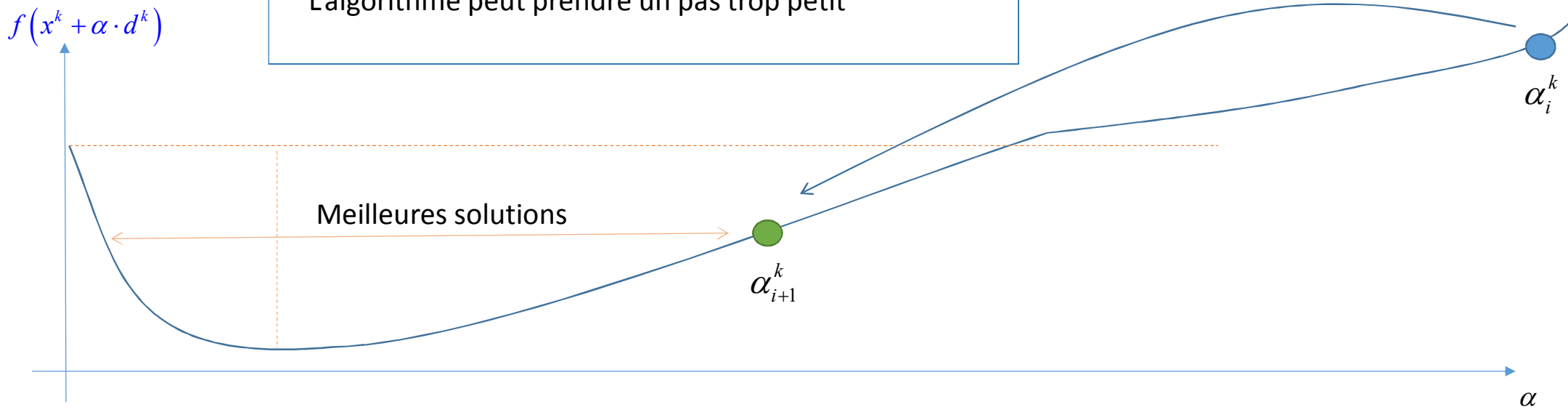
Avantage :

L'algorithme garanti qu'un pas sera trouvé

Inconvénient :

L'algorithme ne permet pas de diminuer  $f$  suffisamment

L'algorithme peut prendre un pas trop petit

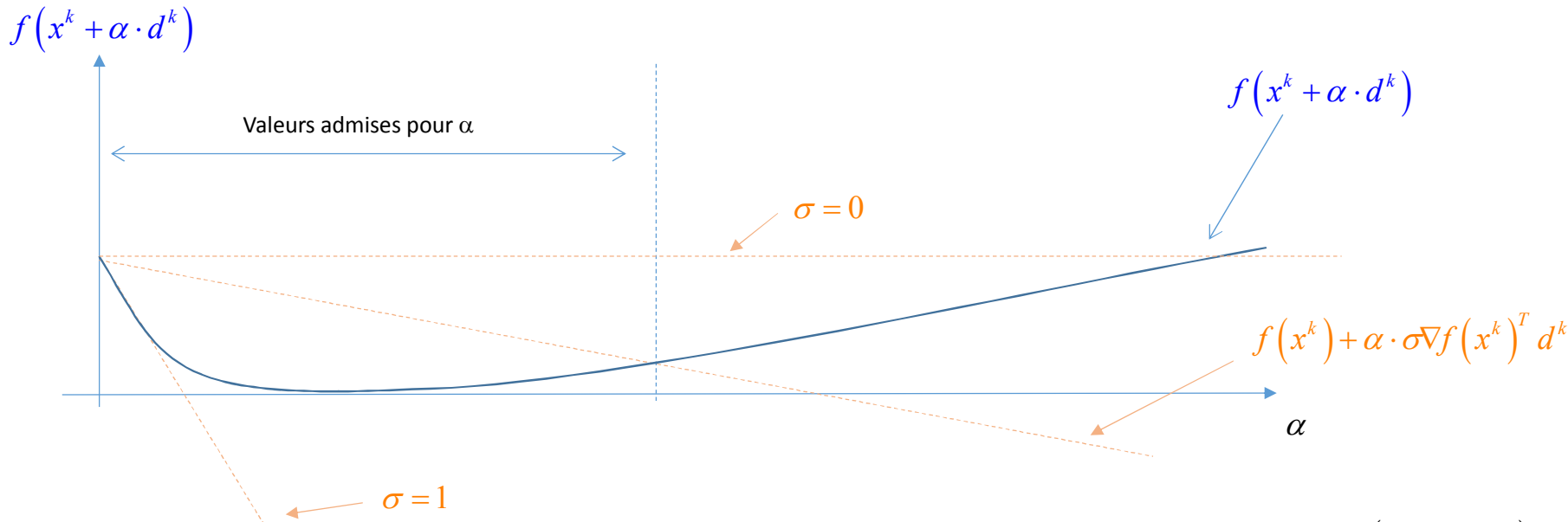


# Algorithme générique : calcul du pas $\alpha^k$ — Règle d'Armijo

- Règle d'Armijo:**

/Armijo 1966/

Idée : on cherche  $\alpha$  tel que  $f(x^k + \alpha \cdot d^k) \leq f(x^k) + \alpha \cdot \sigma \cdot \nabla f(x^k)^T d^k$  Eq. d'une droite avec une pente négative



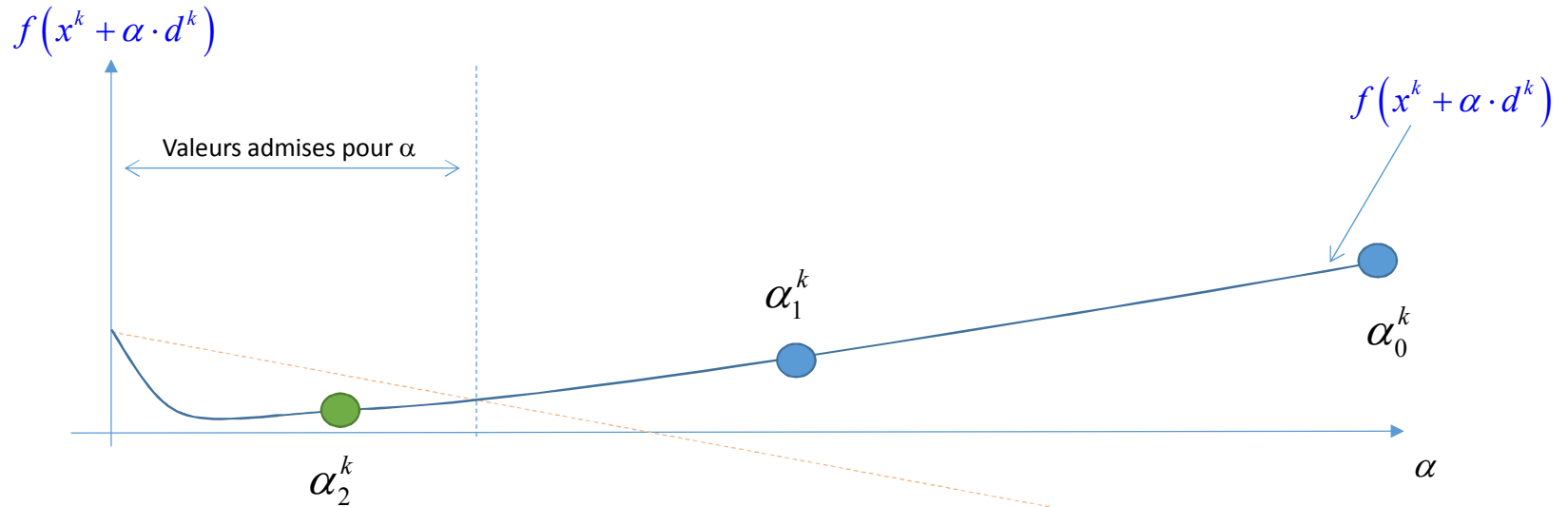
Le paramètre  $\sigma \in ]0, 0.5[$  permet de choisir la pente de la droite, donc la diminution de  $f(x^k + \alpha \cdot d^k) - f(x^k)$  :

- $\sigma \rightarrow 0^+$  Algorithme de *backtracking* mais diminution de  $f(x^{k+1})$  potentiellement faible
- $\sigma \rightarrow 1^-$  Diminution maximale, en général inatteignable car la courbe est souvent au dessus de la tangente
- $\sigma \in ]0, 1[$  Potentiellement « grande » diminution de  $f(x^{k+1})$  mais requiert plus de calculs
- $\sigma \in ]0, 0.5[$  En pratique,  $\sigma < 0.5$  permet de garantir la convergence pour une fonction  $f$  quadratique



## Algorithme générique : calcul du pas $\alpha^k$ — Règle d'Armijo

- Règle d'Armijo:**



### Algorithme au pas $k$ :

Paramètres :  $\alpha_0^k > 0, \sigma \in ]0, 1[$

$i = 0$

Tant que  $f(x^k + \alpha \cdot d^k) \geq f(x^k) + \alpha \cdot \sigma \cdot \nabla f(x^k)^T d^k$

$$\alpha_{i+1}^k = \alpha_i^k \cdot \beta$$

$i = i + 1$

Fin Tant que

$$\alpha^k = \alpha_i^k$$

Avantage: garantie de la décroissance

Inconvénient : le pas peut devenir trop petit

# Algorithme générique : calcul du pas $\alpha^k$ — *Algorithme Backtracking*

Modifications du programme précédent:

```
% Paramètres Armijo
alpha0=1;
beta=0.9;
gamma=0.5;

while ended==0
    dk=-df(x(:,k));

    [alpha(k)]=Armijo(ff,x(:,k),dk,-dk,alpha0,beta,gamma);

    x(:,k+1)=x(:,k)+alpha(k)*dk;
    fxk=fxkp1;
    fxkp1=f(x(:,k+1));

    plot3(x(1,k+1),x(2,k+1),fxkp1,'m*');
    plot3([x(1,k+1) x(1,k)], [x(2,k+1) x(2,k)], [fxkp1 fxk], 'g');

    ended=(k>IterMax) || ...
        norm(fxkp1-fxk)<TolF || ...
        norm(x(:,k)-x(:,k+1))<TolX;
    k=k+1;
end

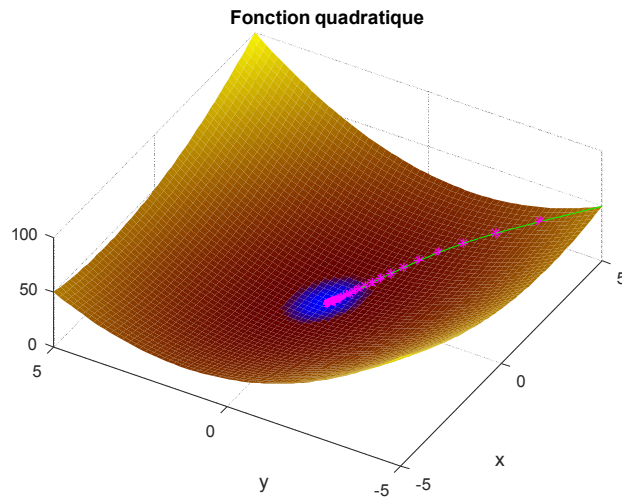
function alpha=Armijo(f,xk,dk,GradF,alpha0,beta,gamma)
alpha=alpha0;
fxk=f(xk);
while f(xk+alpha*dk)>=fxk+alpha*gamma*GradF'*dk
    alpha=alpha*beta;
end
end
```

NB: il est possible d'optimiser le programme pour économiser des appels à la fonction f...

# Exemple académique : fonction quadratique

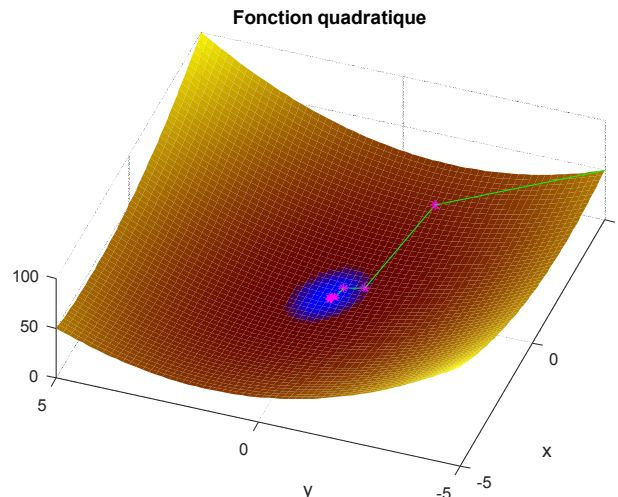
On utilise un pas calculé avec l'algorithme d'Armijo avec:  $\alpha_0^k = 1$   $\sigma = 0.5$

$$\beta = 0.1$$



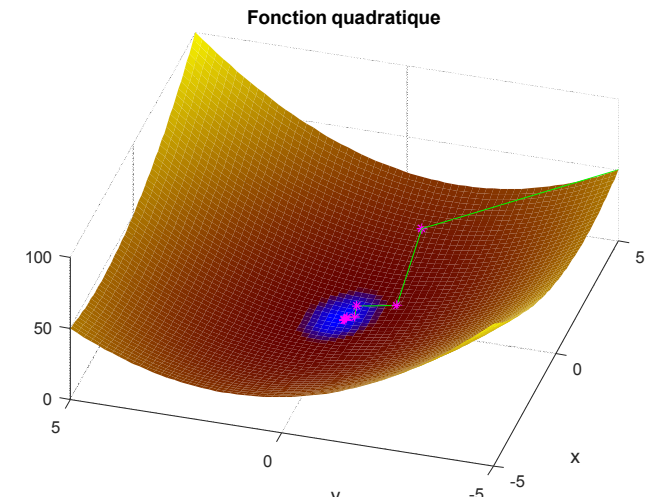
49 itérations boucle principale  
192 appels de f

$$\beta = 0.5$$



11 itérations boucle principale  
45 appels de f

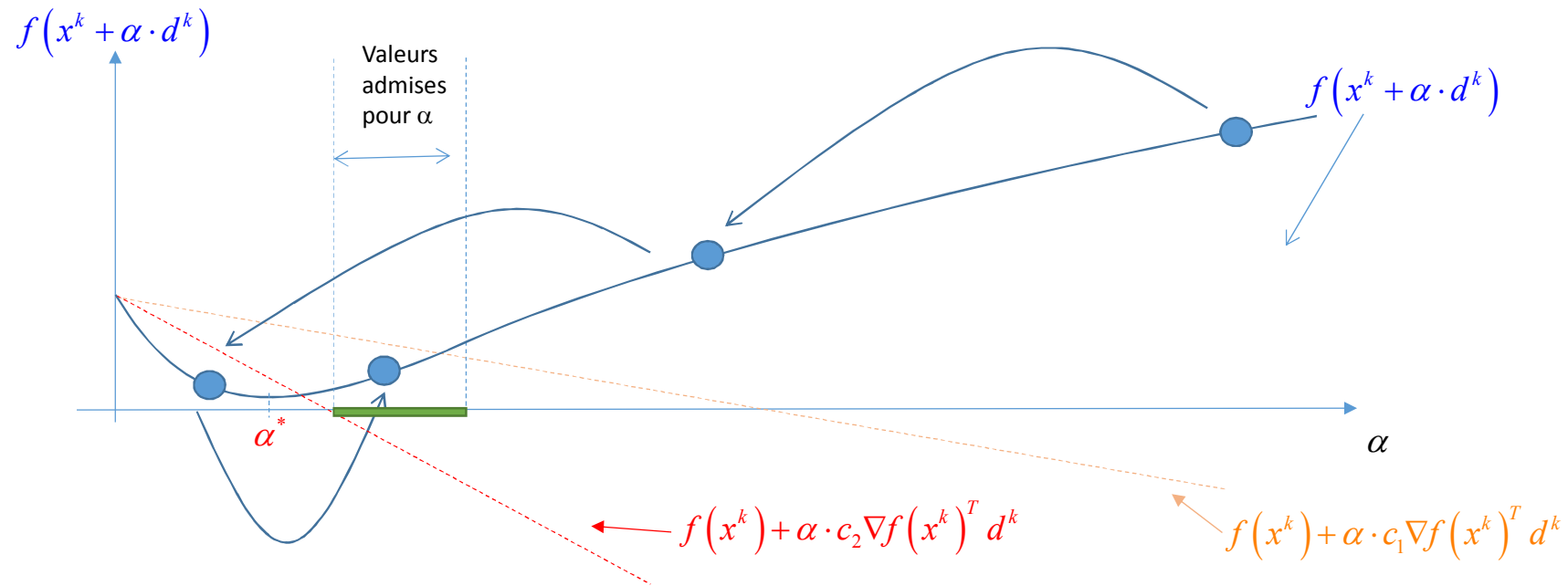
$$\beta = 0.9$$



12 itérations boucle principale  
143 appels de f

Problème: l'algorithme d'Armijo garanti que la fonction  $f$  décroît mais rien ne dit que le pas choisit est proche d'un optimal  
Idée : utiliser les conditions d'optimalité au premier ordre

## Algorithme générique : calcul du pas $\alpha^k$ – *Armijo-Goldstein*



### Règles d'Armijo-Goldstein

Le pas  $\alpha$  doit vérifier:

$$f(x^k + \alpha \cdot d^k) \leq f(x^k) + \alpha \cdot c_1 \cdot \nabla f(x^k)^T d^k$$

$$f(x^k + \alpha \cdot d^k) \geq f(x^k) + \alpha \cdot c_2 \cdot \nabla f(x^k)^T d^k$$

$$0 < c_1 < c_2 < 1$$

Mise en oeuvre: Ne pas accepter les pas trop petits

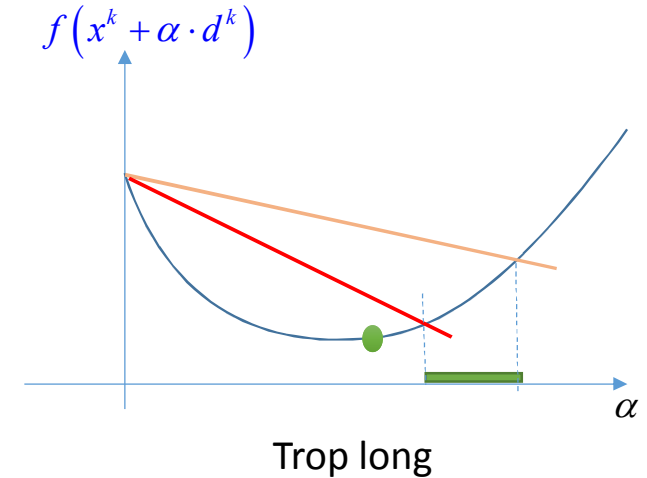
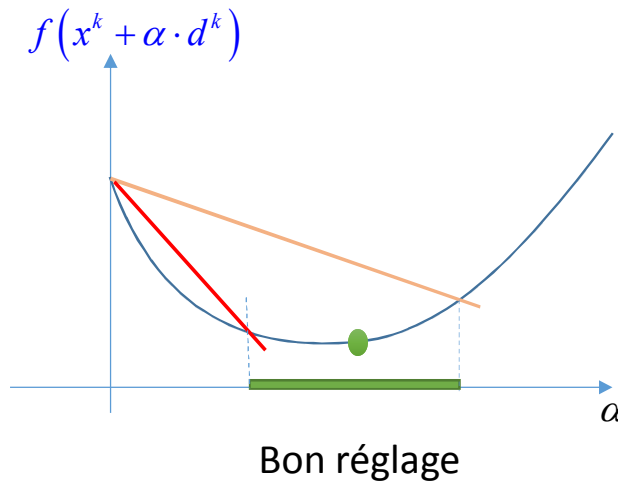
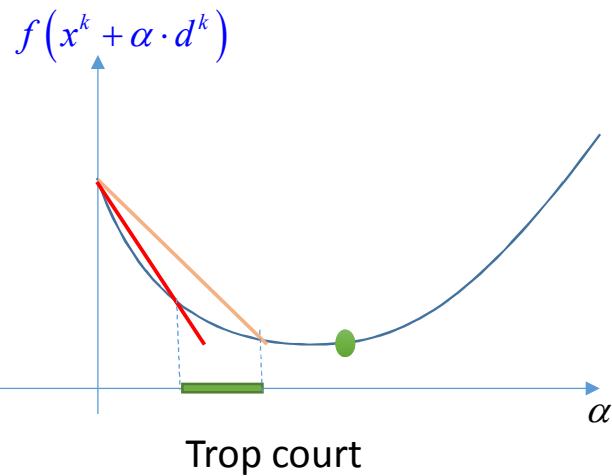
# Algorithme générique : calcul du pas $\alpha^k$ — *Armijo-Goldstein*

```
function alpha=ArmijoGoldstein(f,xk,dk,GradF,alpha0,c1,c2,beta)
alpha=alpha0;
fxk=f(xk);
while (f(xk+alpha*dk)>fxk+alpha*c1*GradF'*dk) || ...
      (f(xk+alpha*dk)<fxk+alpha*c2*GradF'*dk)

    if f(xk+alpha*dk)>=fxk+alpha*c1*GradF'*dk
        alpha=alpha*beta;
    else
        alpha=alpha/beta; % Annule le dernier pas
        beta=beta*0.5;
    end
end
end
```

## Problème lié au réglage des paramètres:

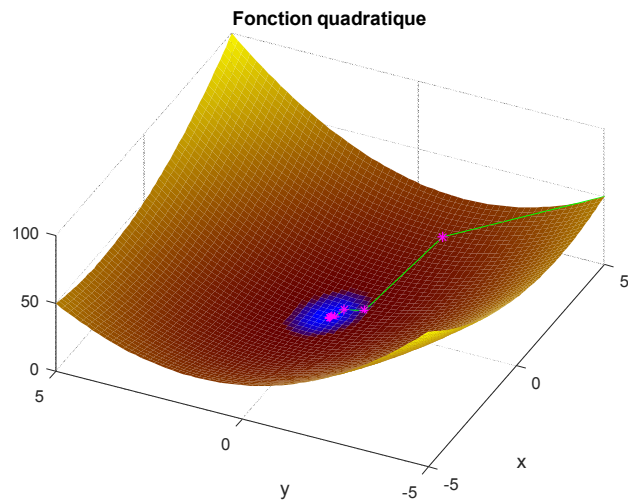
- Suivant les valeurs de  $c_1$  et  $c_2$ , le pas  $\alpha$  peut être trop long ou trop court
- Comment trouver un bon réglage?



## Exemple académique : fonction quadratique

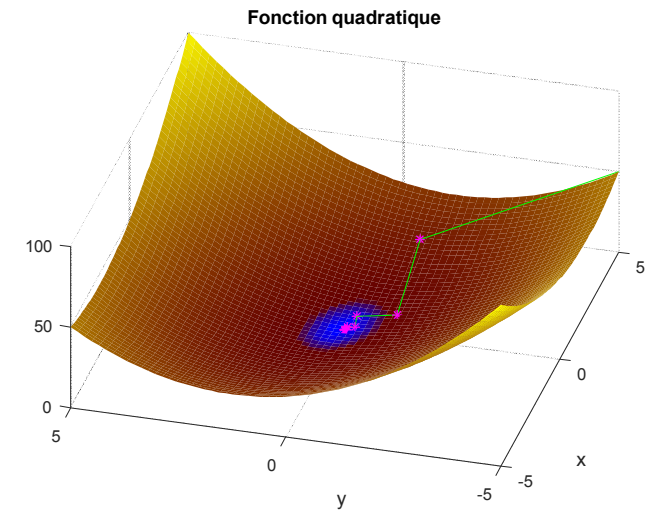
On utilise un pas calculé avec l'algorithme d'Armijo-Goldstein avec:  $\alpha_0^k = 1$   $c_1 = 0.5$   $c_2 = 0.1$

$\beta = 0.5$



11 itérations boucle principale  
60 appels de f

$\beta = 0.9$



9 itérations boucle principale  
82 appels de f

Nb: ne change pas grand-chose sur cet exemple car il n'y avait pas de pb de petits pas

# Algorithme générique : calcul du pas $\alpha^k$ — Règle de Wolfe

Amélioration :

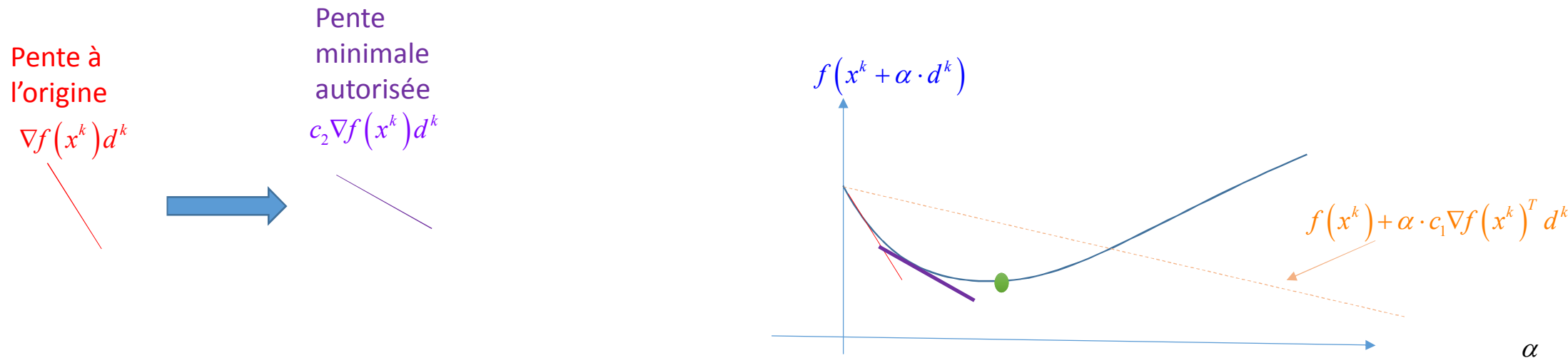
- Tenir compte de la condition d'optimalité au premier ordre :  $\frac{df(x^k + \alpha \cdot d^k)}{d\alpha} = \nabla f(x^k + \alpha \cdot d^k)^T \cdot d^k = 0$

Problème :

- La résolution exacte est trop compliquée => utilisation d'une solution « approché »

Idée :

- Garantir que la borne min permet de rester à gauche de la solution recherchée



# Algorithme générique : calcul du pas $\alpha^k$ — Règle de Wolfe

Règles de wolfe:  $0 < c_1 < c_2 < 1$

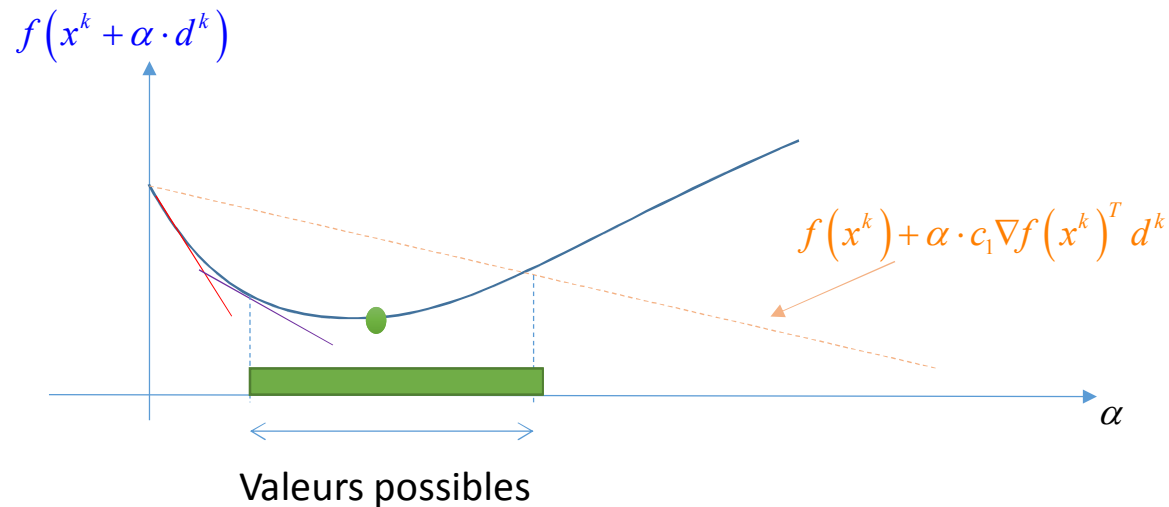
1) La fonction doit décroître suffisamment (Armijo):

$$f(x^k + \alpha \cdot d^k) \leq f(x^k) + \alpha \cdot c_1 \cdot \nabla f(x^k)^T d^k$$

Valeurs typiques :  $c_1=10^{-4}$   $c_2=0,1-0,9$

2) La pente doit être « moins négative » qu'à l'origine

$$\nabla f(x^k + \alpha \cdot d^k)^T d^k \geq c_2 \cdot \nabla f(x^k)^T d^k$$



$$\nabla f(x^k)^T d^k \quad c_2 \nabla f(x^k)^T d^k$$



## Algorithme générique : calcul du pas $\alpha^k$ — Règle de Wolfe

### Règles de wolfe:

$$f(x^k + \alpha \cdot d^k) \leq f(x^k) + \alpha \cdot c_1 \cdot \nabla f(x^k)^T d^k$$

$$\nabla f(x^k + \alpha \cdot d^k)^T d^k \geq c_2 \cdot \nabla f(x^k)^T d^k$$

Valeurs typiques :  $c_1=10^{-4}$   $c_2=0,1-0,9$

⇒ La mise en œuvre est assez difficile

⇒ Trouver un code existant :

Par exemple « *Optimization tutorial* » de Mark Bangert sur Matlab Central

Un autre code : <http://www4.ncsu.edu/~kksivara/ma706>

Un autre code : <https://github.com/GuipengLi>

# Algorithme générique : calcul du pas $\alpha^k$ – Règle de Wolfe

Algorithme proposé par /Jorge Nocedal & Stephen J. Wright, Numerical Optimization/

$$f(x_k + \alpha_k p_k) \leq f(x_k) + c_1 \alpha_k \nabla f_k^T p_k, \quad (3.6a)$$

$$\nabla f(x_k + \alpha_k p_k)^T p_k \geq c_2 \nabla f_k^T p_k, \quad (3.6b)$$

**Algorithm 3.5** (Line Search Algorithm).

Set  $\alpha_0 \leftarrow 0$ , choose  $\alpha_{\max} > 0$  and  $\alpha_1 \in (0, \alpha_{\max})$ ;

$i \leftarrow 1$ ;

**repeat**

Evaluate  $\phi(\alpha_i)$ ;

**if**  $\phi(\alpha_i) > \phi(0) + c_1 \alpha_i \phi'(0)$  or  $[\phi(\alpha_i) \geq \phi(\alpha_{i-1})$  and  $i > 1]$

$\alpha_* \leftarrow \text{zoom}(\alpha_{i-1}, \alpha_i)$  and **stop**;

Evaluate  $\phi'(\alpha_i)$ ;

**if**  $|\phi'(\alpha_i)| \leq -c_2 \phi'(0)$

**set**  $\alpha_* \leftarrow \alpha_i$  and **stop**;

**if**  $\phi'(\alpha_i) \geq 0$

**set**  $\alpha_* \leftarrow \text{zoom}(\alpha_i, \alpha_{i-1})$  and **stop**;

Choose  $\alpha_{i+1} \in (\alpha_i, \alpha_{\max})$ ;

$i \leftarrow i + 1$ ;

**end (repeat)**

**Algorithm 3.6** (zoom).

**repeat**

Interpolate (using quadratic, cubic, or bisection) to find a trial step length  $\alpha_j$  between  $\alpha_{lo}$  and  $\alpha_{hi}$ ;

Evaluate  $\phi(\alpha_j)$ ;

**if**  $\phi(\alpha_j) > \phi(0) + c_1 \alpha_j \phi'(0)$  or  $\phi(\alpha_j) \geq \phi(\alpha_{lo})$

$\alpha_{hi} \leftarrow \alpha_j$ ;

**else**

Evaluate  $\phi'(\alpha_j)$ ;

**if**  $|\phi'(\alpha_j)| \leq -c_2 \phi'(0)$

**Set**  $\alpha_* \leftarrow \alpha_j$  and **stop**;

**if**  $\phi'(\alpha_j)(\alpha_{hi} - \alpha_{lo}) \geq 0$

$\alpha_{hi} \leftarrow \alpha_{lo}$ ;

$\alpha_{lo} \leftarrow \alpha_j$ ;

**end (repeat)**

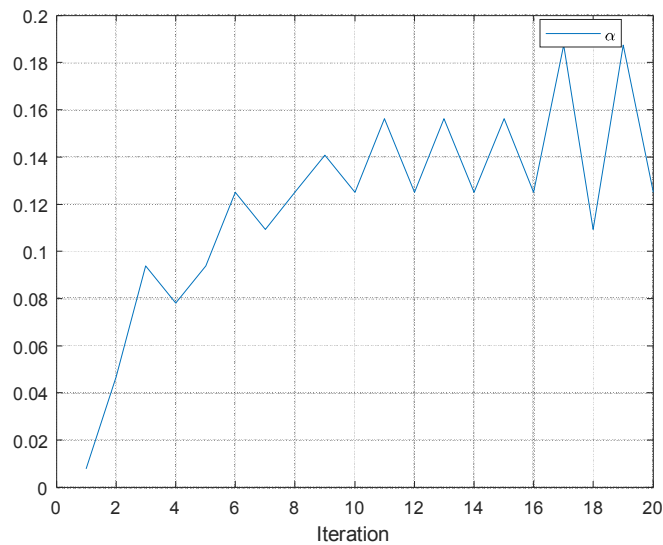
# Exemple académique : fonction quadratique

On utilise un pas calculé avec l'algorithme de Wolfe:

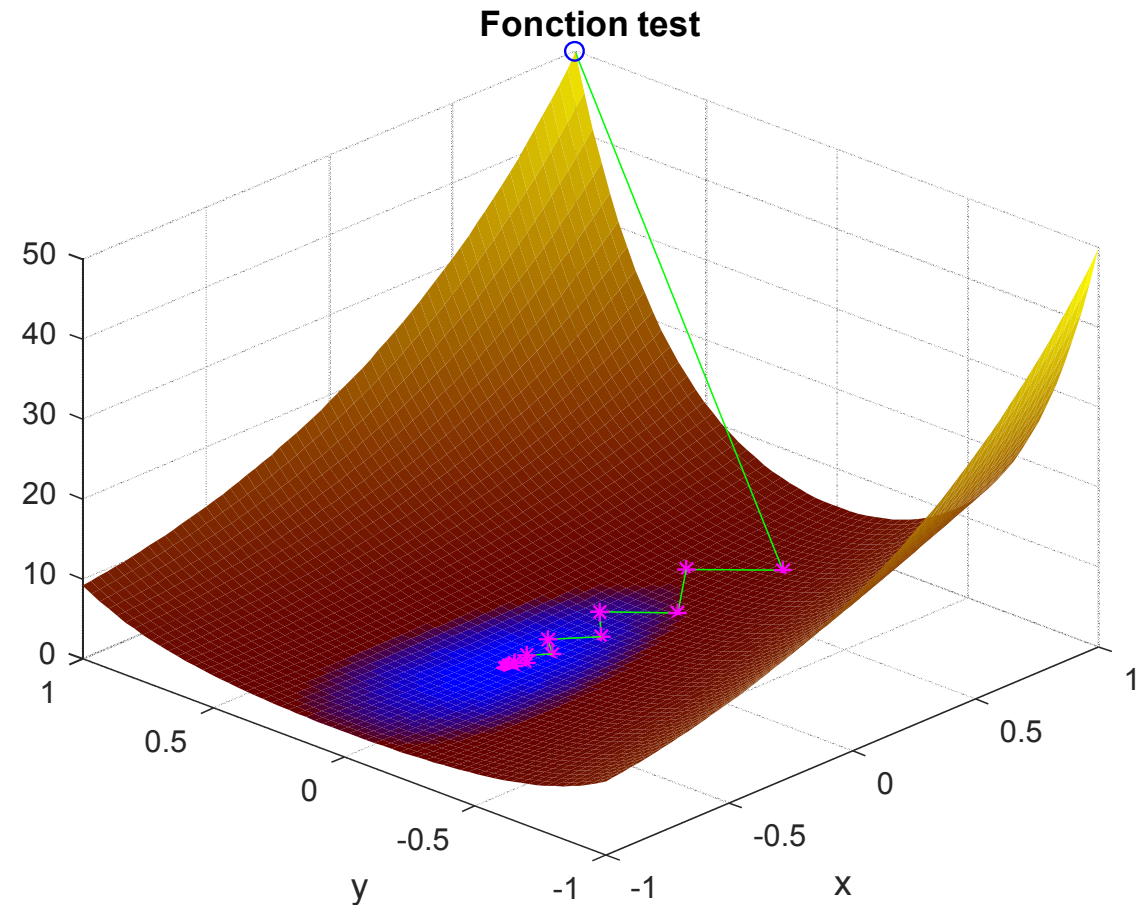
$$\alpha_0^k = 1$$

$$c_1 = 10^{-4}$$

$$c_2 = 0.1$$



21 itérations boucle principale  
151 appels de f



## Chapitre 2 – Optimisation sans contrainte

---

- 1) Contexte
- 2) Algorithme à base de Gradient
- 3) Méthode de Newton & quasi-Newton
- 4) Méthode « Régions de confiance » (Trust region)

# Méthode de Newton

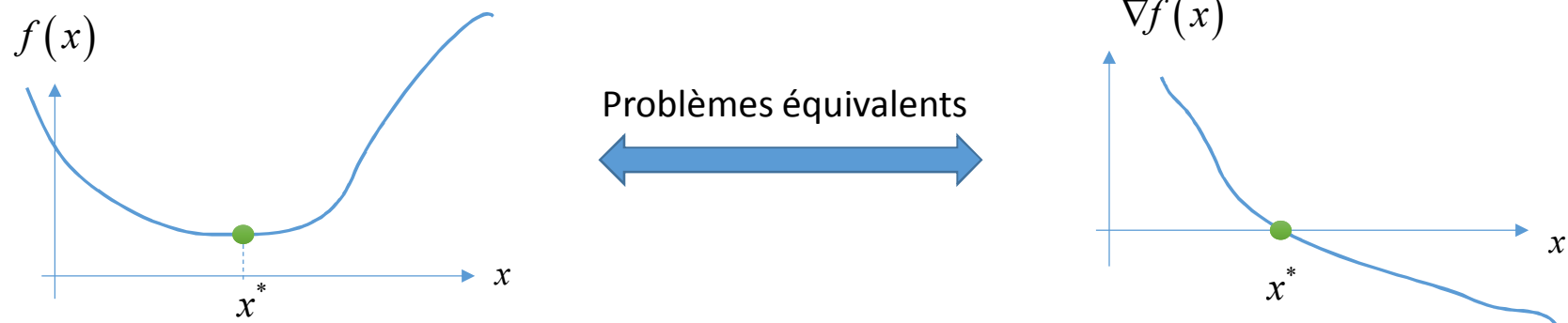
Les algorithmes précédents n'exploitent qu'une information à l'ordre 1.

=> L'algorithme progresse dans une direction qui dépend du gradient

Pour avancer plus vite, une approche consiste à utiliser la condition d'optimalité :

$$\nabla f(x^*) = 0$$

=> Le problème de recherche du minimum est transformé en la recherche d'un zéro d'une fonction multivariables.



# Méthode de Newton

Rappel : Méthode de Newton pour la recherche du zéro d'une fonction

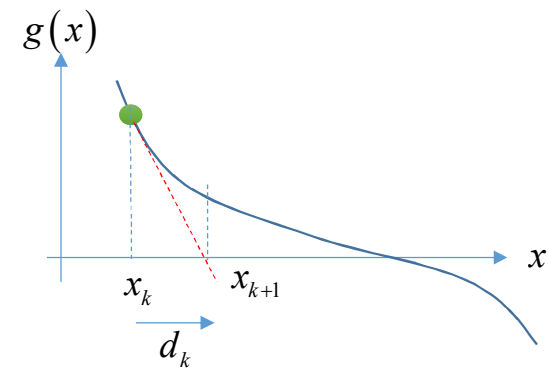
$$g(x^k + d^k) \approx g(x^k) + \nabla g(x^k) d^k$$

$$g(x^k + d^k) = 0 \Leftrightarrow -g(x^k) \approx \nabla g(x^k) d^k$$

D'où la méthode de Newton

Résoudre:  $-g(x^k) = \nabla g(x^k) d^k$

Itérer  $x^{k+1} = x^k + d^k$



Condition d'optimalité au premier ordre pour  $\min_{x \in \mathbb{R}^n} f(x)$

$$\nabla f(x^*) = 0$$

Idée: Appliquer la méthode qui recherche les zéros du gradient  $\nabla f(x^*) = 0$

On pose  $g(x) = \nabla f(x)$  d'où 
$$\begin{cases} -\nabla f(x^k) \approx \nabla^2 f(x^k) d^k \\ x^{k+1} = x^k + d^k \end{cases}$$

# Méthode de Newton

Algorithme de base:

Paramètres :

- $x^0$  : solution initiale
- $\varepsilon$  : tolérance

$k = 0$

Répéter

Chercher  $d^k$  tel que  $-\nabla f(x^k) \approx \nabla^2 f(x^k) d^k$

$$x^{k+1} = x^k + d^k$$

$$k = k + 1$$

Tant que  $\|\nabla f(x^k)\| < \varepsilon$

Hessien

Les problèmes à surmonter:

- Résolution de  $-\nabla f(x^k) \approx \nabla^2 f(x^k) d^k$

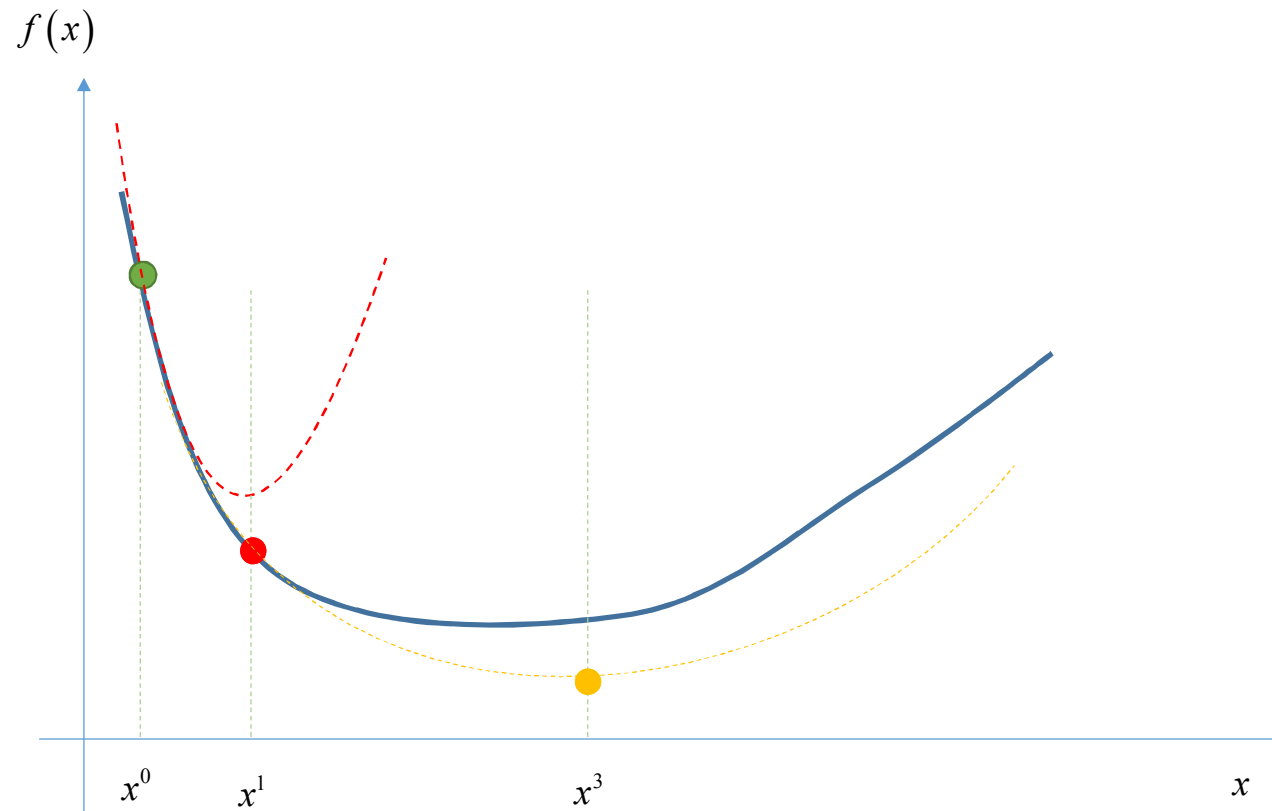
Le Hessien n'est pas nécessairement inversible

- Pas de distinction entre minimum ou maximum local
- Le Hessien n'est pas forcément connu

Système linéaire en  $d^k$

# Méthode de Newton

En fait l'algorithme revient à faire une approximation à l'ordre 2 de  $f$  et à chercher itérativement le minimum.





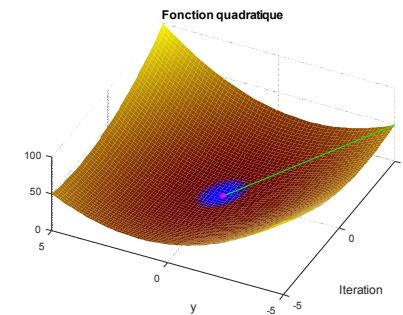
# Exemple académique : fonction quadratique

Sur un exemple quadratique, la méthode fonctionne parfaitement  
(aux erreurs numériques près)

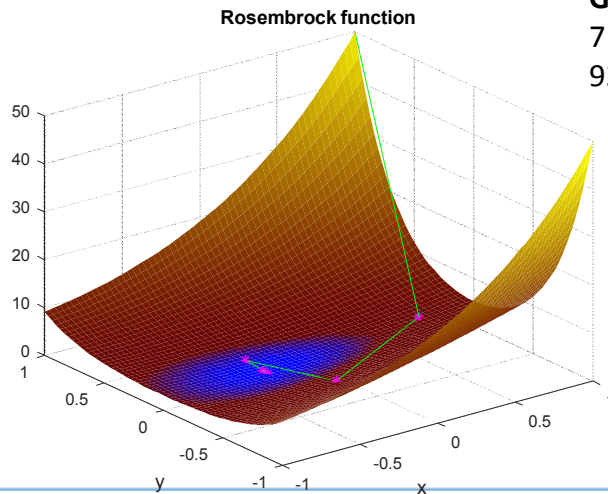
=> Nécessité de comparer sur un exemple plus complexe

En général, fonctionne mal car des problèmes numériques  
surviennent lors de la résolution du système linéaire

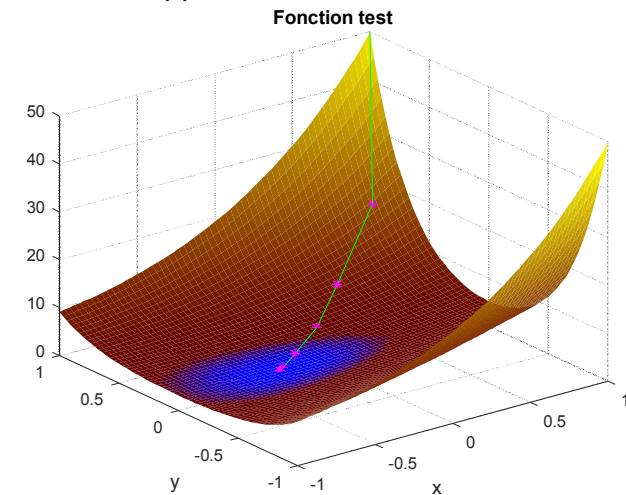
$$f(x,y) = e^{x-3y-\frac{1}{10}} + e^{x+3y-\frac{1}{10}} + e^{-x-\frac{1}{10}}$$



3 itérations boucle principale  
2 appels de f



**GRADIENT+ARMIJO**  
7 itérations boucle principale  
92 appels de f



**NEWTON**  
9 itérations boucle principale  
9 appels de f

# Méthode de quasi-Newton

Idée:

Si le Hessien est inversible  $-\nabla f(x^k) \approx \nabla^2 f(x^k) d^k \Leftrightarrow d^k = -(\nabla^2 f(x^k))^{-1} \cdot \nabla f(x^k)$   
 $x^{k+1} = x^k + d^k$

On retrouve une forme « similaire » à celle de l'algorithme à base de gradient  
 $d^k$  est une direction de descente si et seulement si le Hessien est défini positif

On remplace l'inverse du Hessien, souvent inconnue, par une approximation:

$$d^k = -\alpha^k S^k \cdot \nabla f(x^k)$$
$$x^{k+1} = x^k + d^k$$

=> Les algorithmes doivent « rendre » l'approximation du Hessien positive

Différentes approximations donnent différentes méthodes:

- Davidson-Fletcher-Powell (DFP).
- Levenberg-Marquardt
- Broyden-Fletcher-Goldfarb-Shanno (BFGS)

# Méthode de quasi-Newton : *algorithme de Davidson-Fletcher-Powell (DFP)*

Idée: On estime non pas le Hessien, mais directement son inverse

Paramètres:

- $S^0$  : matrice définie positive

$$k = 0$$

Tant que **Critère d'arrêt** non satisfait

$$d^k = -S^k \nabla f(x^k)$$

Choisir  $\alpha^k$  par une méthode de calcul de pas (Armijo, etc.)

$$x^{k+1} = x^k + \alpha^k d^k$$

$$\gamma^k = \nabla f(x^{k+1}) - \nabla f(x^k)$$

$$\delta^k = \alpha^k d^k$$

$$S^{k+1} = S^k + \frac{\delta^k (\delta^k)^T}{(\delta^k)^T \gamma^k} - \frac{S^k \gamma^k (S^k \gamma^k)^T}{(\gamma^k)^T S^k \gamma^k}$$

# Méthode de quasi-Newton : *algorithme de Davidson-Fletcher-Powell (DFP)*

```

IterAlpha=0;
Sk=eye(size(x,1));
while ended==0
    dfk=df(x(:,k));
    dk=-Sk*dfk;
    alpha=Armijo(f2,x(:,k),dk,dfk,alpha0,beta,gamma);

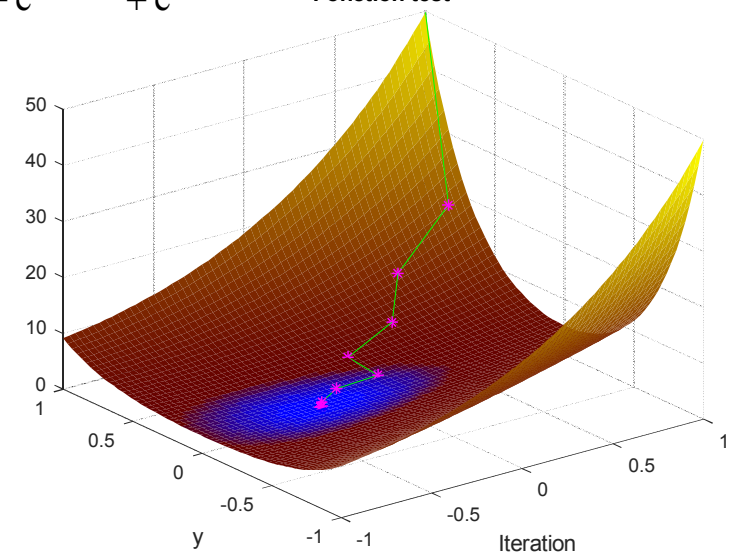
    x(:,k+1)=x(:,k)+alpha*dk;
    Gammak=df(x(:,k+1))-dfk;
    deltak=alpha*dk;
    Sk=Sk+(deltak*deltak')/(deltak'*Gammak)-(Sk*Gammak)*(Sk*Gammak')/(Gammak'*Sk*Gammak);

    ended=(k>IterMax) || ...
        norm(df(x(:,k)))<EpsGrad || ...
        norm(x(:,k)-x(:,k+1))<Tolx;
    k=k+1;
end

```

$$f(x,y) = e^{x-3y-\frac{1}{10}} + e^{x+3y-\frac{1}{10}} + e^{-x-\frac{1}{10}}$$

Fonction test



**DFP+ARMIJO**  
 14 itérations boucle principale  
 52 appels de f

# Méthode de quasi-Newton : *Broyden-Fletcher-Goldfarb-Shanno (BFGS)*

Paramètres:

- $S^0$  : matrice définie positive

$$k = 0$$

Tant que **Critère d'arrêt** non satisfait

$$d^k = -S^k \nabla f(x^k)$$

Choisir  $\alpha^k$  par une méthode de calcul de pas (Armijo, etc.)

$$x^{k+1} = x^k + \alpha^k d^k$$

$$\gamma^k = \nabla f(x^{k+1}) - \nabla f(x^k)$$

$$\delta^k = \alpha^k d^k$$

$$S^{k+1} = \left( I - \frac{\delta^k (\gamma^k)^T}{(\delta^k)^T \gamma^k} \right) S^k \left( I - \frac{\delta^k (\gamma^k)^T}{(\delta^k)^T \gamma^k} \right) + \frac{\delta^k (\delta^k)^T}{(\delta^k)^T \gamma^k}$$

# Méthode de quasi-Newton : *Broyden-Fletcher-Goldfarb-Shanno (BFGS)*

```

IterAlpha=0;
n=size(x,1);
Sk=eye(n);
while ended==0
    dfk=df(x(:,k));
    dk=-Sk*dfk;
    alpha=Armijo(f2,x(:,k),dk,dfk,alpha0,beta,gamma);

    x(:,k+1)=x(:,k)+alpha*dk;
    Gammak=df(x(:,k+1))-dfk;
    deltak=alpha*dk;
    Sk=(eye(n)-deltak*Gammak'/(deltak'*Gammak))*Sk*(eye(n)-deltak*Gammak'/(deltak'*Gammak))+deltak*deltak'/(deltak'*Gammak);
    fxk=fxkp1;
    fxkp1=f(x(:,k+1));

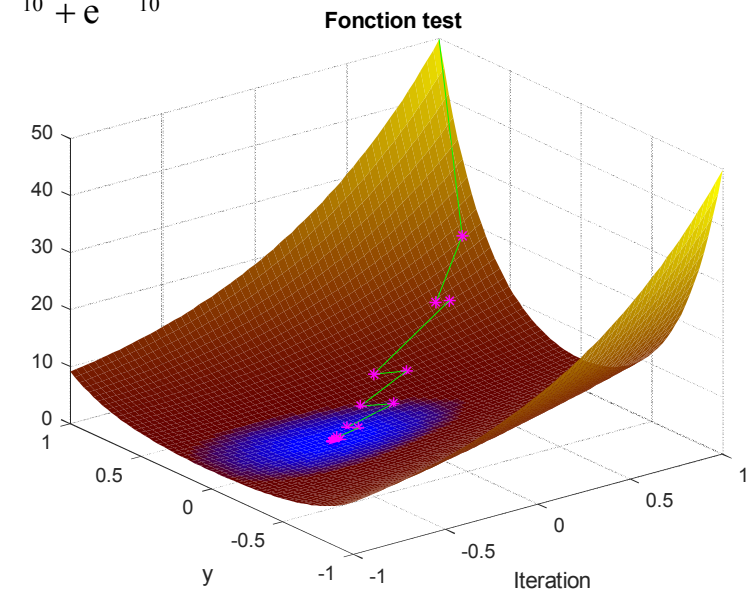
    ended=(k>IterMax)||...
        norm(df(x(:,k)))<EpsGrad||...
        norm(x(:,k)-x(:,k+1))<Tolx;
    k=k+1;
end

```

$$f(x,y)=e^{x-3y-\frac{1}{10}}+e^{x+3y-\frac{1}{10}}+e^{-x-\frac{1}{10}}$$

## **BFGS+ARMIJO**

22 itérations boucle principale  
105 appels de f

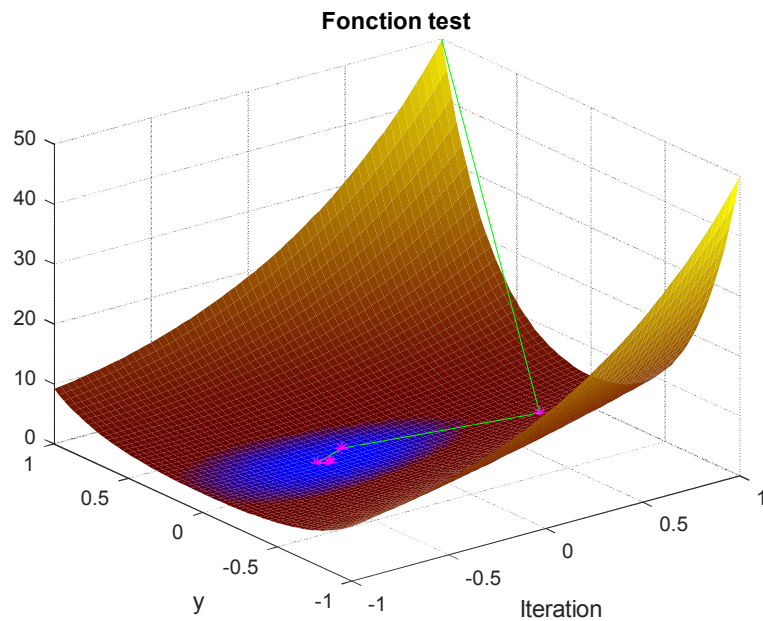


# Méthode de quasi-Newton : *Broyden-Fletcher-Goldfarb-Shanno (BFGS)*

En pratique, l'algorithme BFGS pour la direction + Règles de Wolfe pour le choix du pas est souvent très efficace

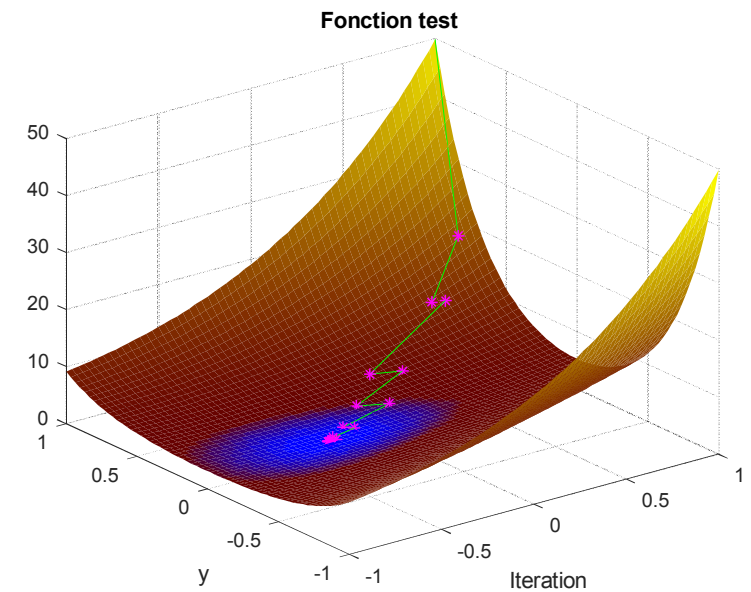
## **BFGS+Wolfe**

10 itérations boucle principale  
56 appels de f



## **BFGS+ARMIJO**

22 itérations boucle principale  
105 appels de f



# Résumé

Algorithme d'optimisations:  $x^{k+1} = x^k + \alpha^k d^k$

- Trouver un pas  $\alpha^k$
- Trouver une direction  $d^k$

Algorithme de choix du pas	Performances	Facilité de mise en œuvre
Pas optimal	++	Impossible sauf sur quelques cas simples
Pas constant	----	++++
Backtracking	+	+++
Armijo	++	++
Armijo-Golstein	++	+
Wolfe	++++	--

Algorithme de choix de la direction	Performances	Facilité de mise en œuvre
Plus forte pente (Gradient)	+	+++ Nécessite le gradient
Newton	++	Nécessite un Hessien défini positif
DFP	+++	+
BFGS	++	+++



## Chapitre 2 – Optimisation sans contrainte

---

- 1) Contexte
- 2) Algorithme à base de Gradient
- 3) Méthode de Newton & quasi-Newton
- 4) Méthode « Régions de confiance » (Trust region)

Référence utilisée pour ce cours: Tamas Terlaky, Advanced Optimization Lab., CAS, McMaster 9

# Régions de confiance

Autour de  $x^k$ , la fonction  $f$  peut être approché (localement) par une quadratique

$$q(x) \approx f(x^k) + \nabla f(x^k)^T (x - x^k) + \frac{1}{2} (x - x^k)^T \nabla^2 f(x^k) (x - x^k)$$

On suppose que l'approximation est valide pour une région dite de confiance:

$$\|x - x^k\| < \delta^k$$

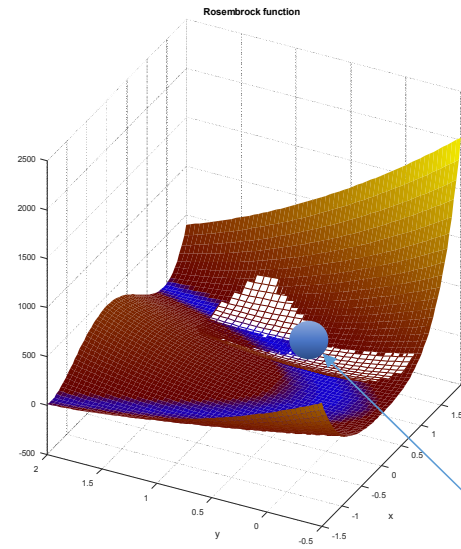
La forme dépend de la norme choisie :

- $\|\cdot\|_2$ : cercle
- $\|\cdot\|_\infty$ : carré

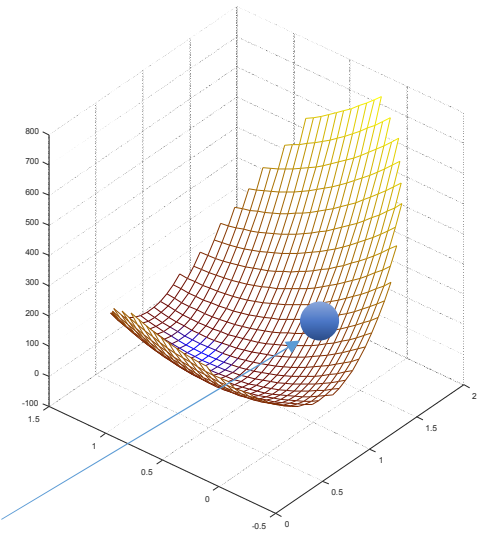
Dans le cours, on utilise  $\|\cdot\|_2$

Exemple : fonction de Rosembrock

$$f(x, y) = (1 - x)^2 + 100(y - x^2)^2$$



Fonction



Modèle local quadratique

# Régions de confiance

## Algorithme de base :

Init :  $k = 0$     $x^0 = x_0$     $\delta^0 = 1$

Tant que Critère non satisfait

$$s^k = \arg \min_{\|s^k\| \leq \delta^k} q(x^k + s^k) = f(x^k) + \nabla f(x^k)^T s^k + \frac{1}{2} (s^k)^T \nabla^2 f(x^k) s^k$$

$$Ratio = \frac{f(x^k + s^k) - f(x^k)}{q(x^k + s^k) - q(x^k)}$$

Variation de la fonction  
Variation du modèle

Mesure la confiance que l'on peut avoir dans le modèle

Si  $Ratio < \eta_1$  alors

$\delta^k = \delta^k \cdot \gamma_{red}$  *Echec, on réduit la taille de la région*

Sinon

$x^{k+1} = x^k + s^k$  *Réussite, on accepte la solution obtenue*

Si  $Ratio > \eta_2$  alors

$\delta^{k+1} = \delta^k \cdot \gamma_{aug}$  *« Grande » Réussite, on augmente la taille de la région de confiance*

Sinon

$\delta^{k+1} = \delta^k$  *Réussite « modeste »: On garde la même taille de région*

Fin si

$k = k + 1$

Fin si

Fin tant que

Paramètres    $\eta_1 < \eta_2 < 1$   
 $\gamma_{red} < 1$     $\gamma_{aug} > 1$

# Régions de confiance

```
clear all;
close all;
clc;

NoFonction=3;
ChoixFonction;

x0=[0; 0.5];

% Affichage de la fonction
X=linspace(xmin,xmax,61);
Y=linspace(ymin,ymax,62);
Z=zeros(length(Y),length(X));
for i=1:length(X)
    for j=1:length(Y)
        Z(j,i)=f([X(i);Y(j)]);
    end
end

figure;
surf(X,Y,Z,'EdgeColor','none');
hold on;
xlabel('x');
ylabel('y');
title(Fname)
%colorbar;
colormap(map)
hold on;
plot3(x0(1),x0(2),f(x0),'bo');

% Criteres d'arret
Tolf=1e-6;
Tolx=1e-9;
IterMax=100;
```

```
while ended==0
    xk=x(:,k); % sol courante
    fxk=f(xk); % valeur de la fonction
    g=df(xk); % gradient
    H=d2f(xk); % hessien

    % Modèle au point xk
    q=@(x2) fxk+g'(x2-x(:,k))+0.5*(x2-x(:,k))'*H*(x2-x(:,k));

    % Résolution du pb Trust Region
    s(:,k)=xxxxxxxxxxxxxxxxxxxxx ←
    xkp1=xk+s(:,k); % prochain point ?

    plot3(xkp1(1),xkp1(2),f2(xkp1),'r*');
    plot3([xkp1(1) xk(1)],[xkp1(2) xk(2)],[f2(xkp1) fxk],'g');

    % Amélioration sur le Modèle et la fonction
    ModeleDiff=q(xk+s(:,k))-q(xk);
    FonctionDiff=f(xkp1)-fxk;
    Ratio=FonctionDiff/ModeleDiff;

    % Gestion de la taille de la région
    if Ratio<nul
        % Echec
        Delta(k)=Delta(k)*Gamma1;
    else
        % Réussite
        x(:,k+1)=xkp1;
        if Ratio>nu2
            % Grande réussite
            Delta(k+1)=Delta(k)*Gamma2;
        else
            % Réussite modeste
            Delta(k+1)=Delta(k);
        end
    end
end
```

```
fxkp1=f(x(:,k+1));
ended=(k>IterMax)||...
        norm(fxkp1-fxk)<Tolf||...
        norm(x(:,k)-x(:,k+1))<Tolx;
k=k+1;

end
end
```

Appeler ici la fonction qui calcule s

# Régions de confiance

## Résolution de l'optimisation du problème contraint:

Intérêt de la contrainte: il existe toujours une solution

Nécessite de savoir gérer les problèmes avec contraintes (cf prochain cours)

Le problème d'optimisation :

$$\min J = f(x^k) + \nabla f(x^k)^T s + \frac{1}{2} s^T \nabla^2 f(x^k) s$$

sous la contrainte :  $\|s\| \leq \delta$

A pour solution:  $s(\mu) = -(\nabla^2 f(x^k) + \mu I)^{-1} (\nabla f(x^k))^T$

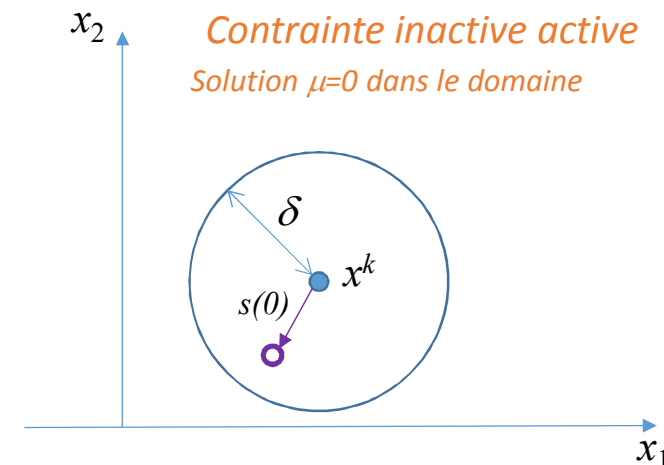
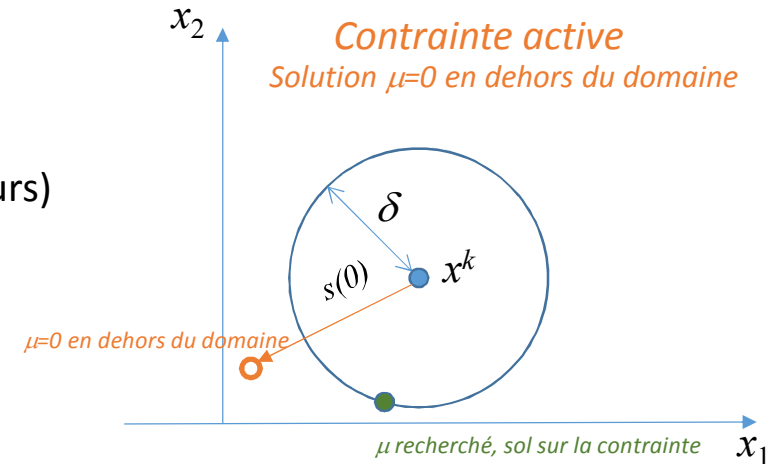
Si  $\|s(0)\| > \delta$  alors

$\mu$  est tel que  $\|s(\mu)\| = \delta$  *Contrainte active*

Sinon

$\mu=0$  *Contrainte non active*

Fsi



# Régions de confiance

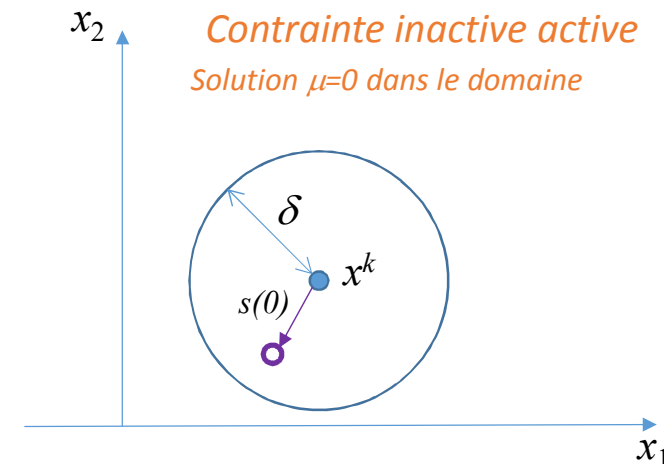
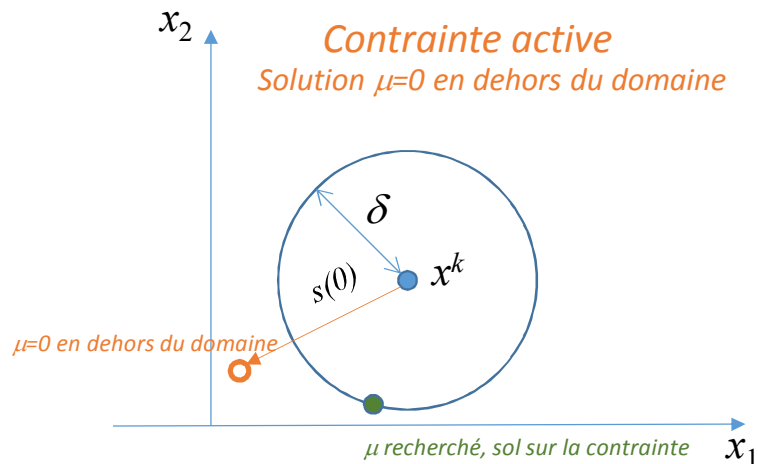
Pourquoi les algorithmes à base de région de confiance sont performants ?

La direction de recherche change avec la contrainte: quand la taille de la région change, la direction de recherche change (elle évolue entre la direction donnée par la méthode de descente maximale et celle de Newton)

=> Permet d'éviter un point de selle local par exemple

$$s(\mu) = -(\nabla^2 f(x^k) + \mu I)^{-1} (\nabla f(x^k))^T$$

$\mu = 0 \rightarrow s(0) = -\nabla^2 f(x^k)^{-1} (\nabla f(x^k))^T \Rightarrow \text{Newton}$ 
 $\mu \rightarrow +\infty \rightarrow s(\mu) \rightarrow -(\nabla f(x^k))^T \Rightarrow \text{Plus forte pente}$



## Régions de confiance : Méthode de Cauchy (plus forte pente)

Méthode de Cauchy:

- 1) On avance en suivant la plus forte pente
- 2) On s'arrête sur la contrainte

Problème non contraint:  $s^{uc} = \arg \min_s q(x^k + s)$

$$q(x^k + s) = f(x^k) + g^T s + \frac{1}{2} s^T H s \quad H = \nabla^2 f(x^k) \quad g = \nabla f(x^k)$$

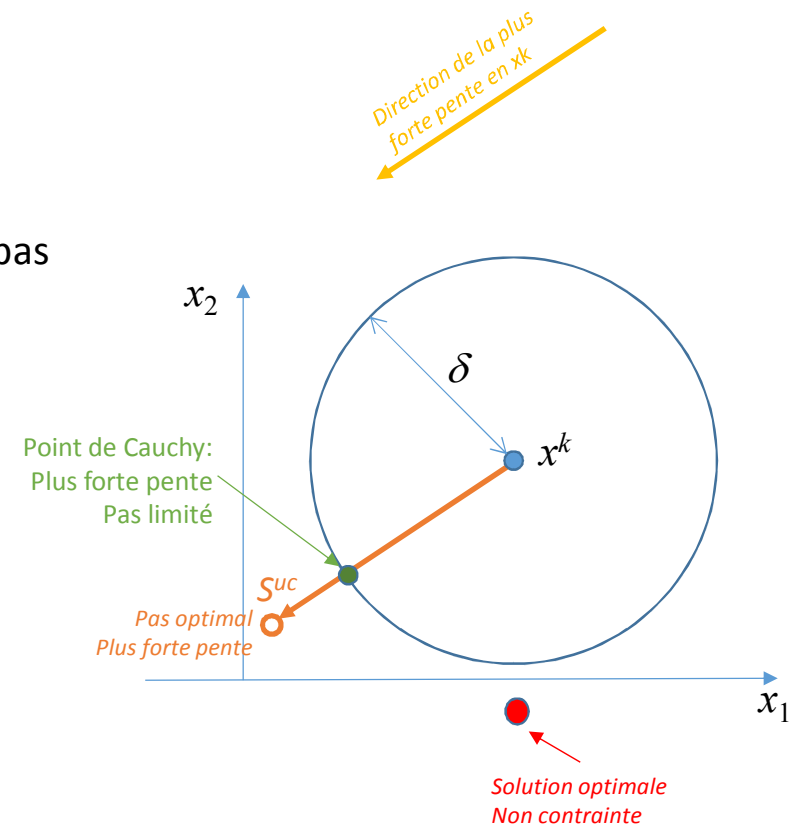
Direction de la plus forte pente:  $s = -t \cdot g$  Avec  $t$  la longueur du pas

$$\Rightarrow q(x^k + s) = f(x^k) - t g^T g + \frac{t^2}{2} g^T H g$$

Choix du pas  $t$  optimal:

$$\frac{\partial q(x^k + s)}{\partial t} = -g^T g + t g^T H g = 0 \quad t = \frac{g^T g}{g^T H g}$$

$$\Leftrightarrow s^{uc} = -\frac{g^T g}{g^T H g} g$$



## Régions de confiance : Méthode de Cauchy (plus forte pente)

### Méthode de Cauchy:

- 1) On avance en suivant la plus forte pente
- 2) On s'arrête sur la contrainte

### Algorithme

Si  $g^T Hg > 0$  alors

*La fonction est localement convexe*

$$t = \frac{g^T g}{g^T Hg} \quad s^{uc} = -t \cdot g$$

$$\text{longueur} = \|s^{uc}\| = \sqrt{(s^{uc})^T (s^{uc})} = t \cdot \|g\|$$

Si longueur  $< \delta$  alors

$$s = s^{uc} \quad \text{Sol non contrainte est ok}$$

Sinon

$$t = \delta$$

$$s = -t \cdot \frac{g}{\|g\|} \quad \text{On doit limiter la longueur du pas}$$

Fin si

Sinon

*La fonction n'est pas localement convexe*

$$t = \delta$$

$$s = -t \cdot \frac{g}{\|g\|} \quad \text{On prends le plus grand pas possible}$$

Fin si

```
function s=TR_Cauchy(g,H,Delta)
    gHg = g'*H*g;

    if gHg>5*eps,
        % La fonction est localement convexe
        % Choix de la plus forte descente
        t=(g'*g)/gHg;
        s_uc = -t*g;
        longueur = norm(s_uc); % Norme 2
        if longueur <= Delta,
            % Solution retenue
            s=s_uc;
        else
            % On sature la longueur du pas
            s= -Delta*g/norm(g);
        end
    else
        % Modèle localement non convexe:
        % SD minimizer lies on TR boundary, in direction of -g'
        s= -Delta*g/norm(g); % On avance dans la direction de g, longueur : delta
    end
end
```



## Régions de confiance : *Méthode de Cauchy (plus forte pente)*

Exemple:  $f(x, y) = e^{x-3y-\frac{1}{10}} + e^{x+3y-\frac{1}{10}} + e^{-x-\frac{1}{10}}$

### Trust Region + Cauchy

17 itérations boucle principale

49 appels de f

### Intérêt :

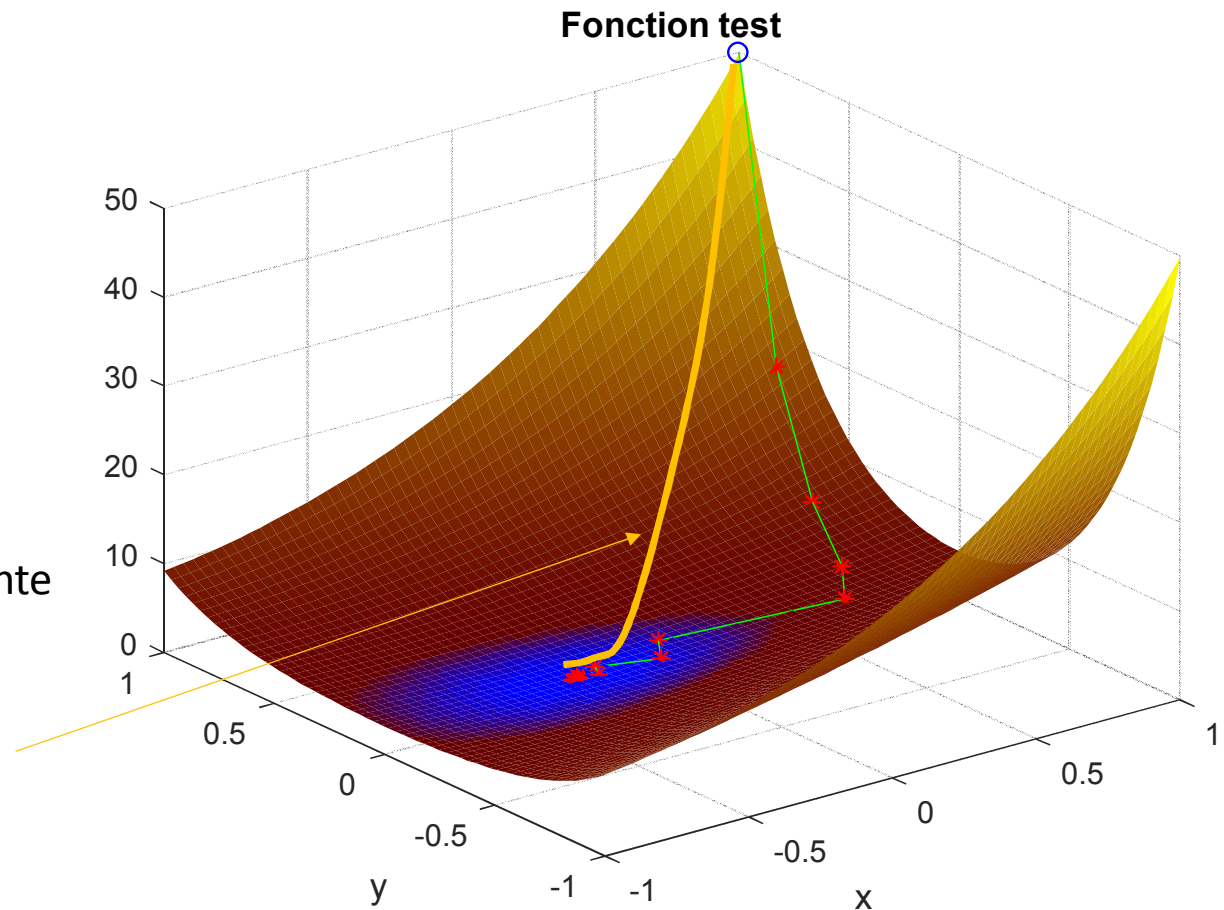
Méthode très robuste en général

### Inconvénient :

On reste sur la direction de la plus forte pente

⇒ Peu mieux faire

Une descente vers le minimum serait  
potentiellement plus efficace

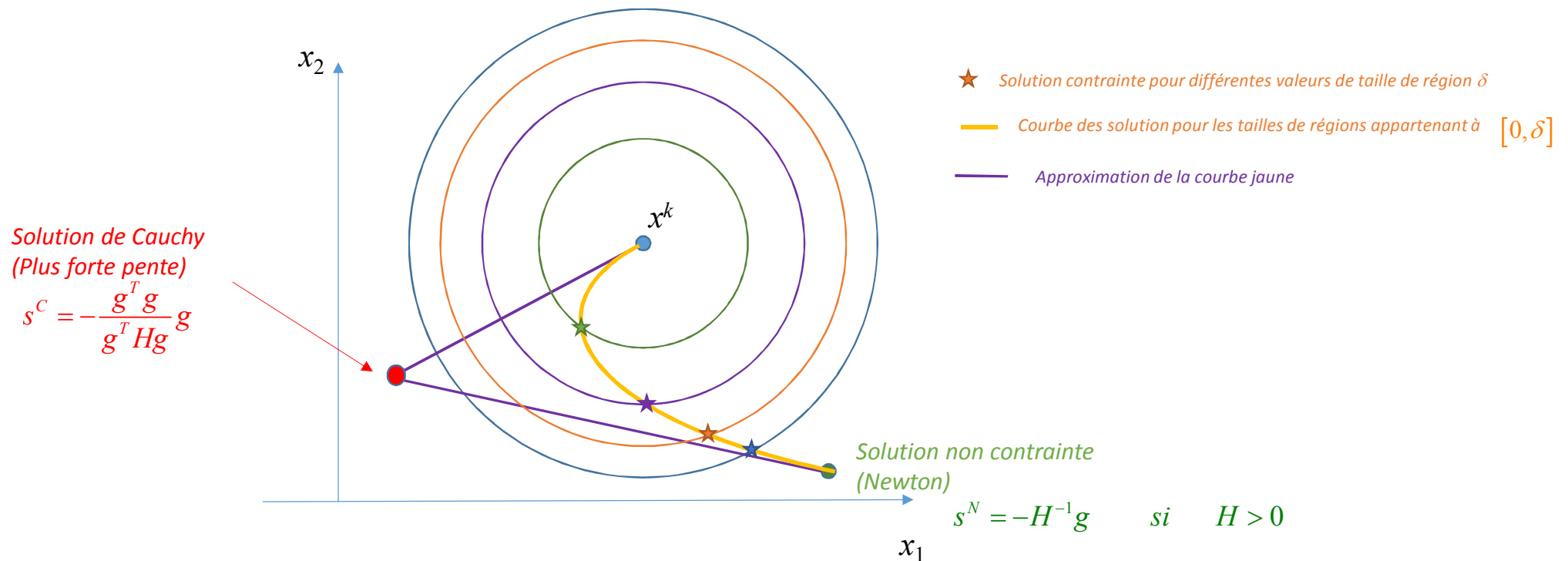


## Régions de confiance : Méthode « Dogleg »

$$\min J = f(x^k) + g^T s + \frac{1}{2} s^T H s \quad H = \nabla^2 f(x^k) \quad g = \nabla f(x^k)$$

sous la contrainte :  $\|s\| \leq \delta$

PB : Calculer l'intersection de la courbe jaune avec la limite de la région est impossible



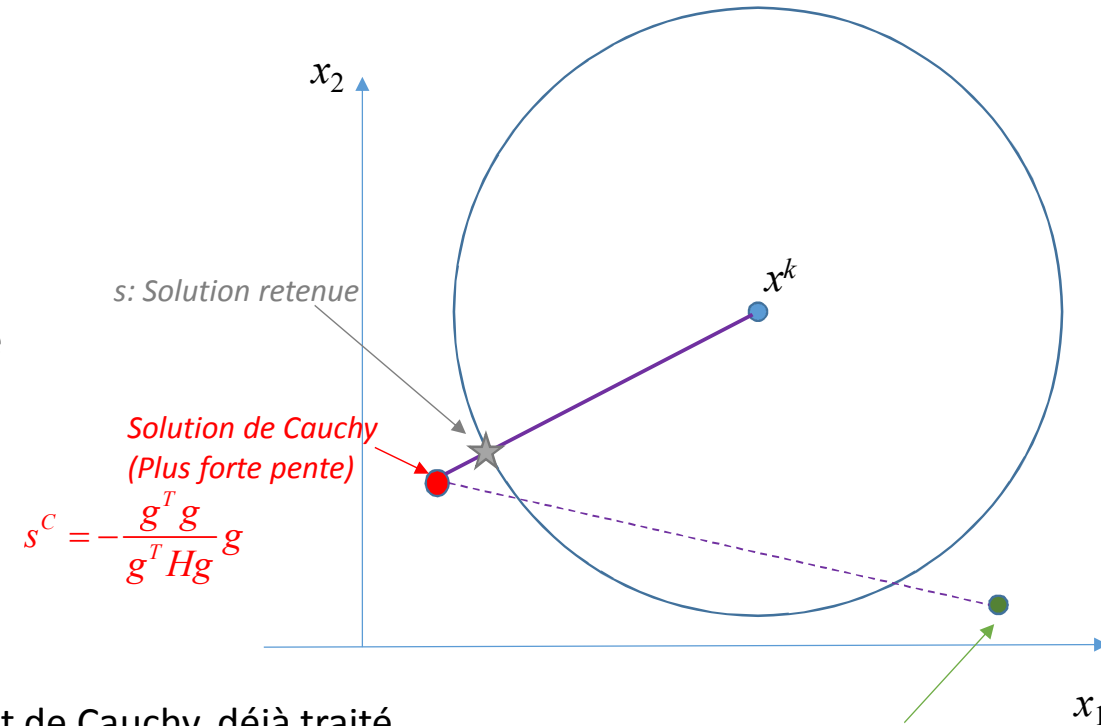
## Régions de confiance : Méthode « Dogleg »

$$\min J = f(x^k) + g^T s + \frac{1}{2} s^T H s \quad H = \nabla^2 f(x^k) \quad g = \nabla f(x^k)$$

sous la contrainte :  $\|s\| \leq \delta$

On approche la courbe par 2 segments de droites

=> Calcul de l'intersection entre une droite et un cercle facile



Cas n° 1 : Intersection entre  $[x^k, s^C]$  et le cercle : C'est le point de Cauchy, déjà traité

$$s = -\delta \cdot \frac{g}{\|g\|}$$

Direction du gradient (normalisée)

Pas maximal autorisé:  $\delta$

*Solution non contrainte (Newton)*

$$s^N = -H^{-1} g \quad \text{si} \quad H > 0$$

## Régions de confiance : Méthode « Dogleg »

$$\min J = f(x^k) + g^T s + \frac{1}{2} s^T H s \quad H = \nabla^2 f(x^k) \quad g = \nabla f(x^k)$$

sous la contrainte :  $\|s\| \leq \delta$

Cas n° 2 : Intersection entre le segment  $[s^N, s^c]$  et le cercle

Equation d'un point sur le segment  $[s^N, s^c]$ :

$$p(t) = x^k + s^c + t(s^N - s^c) \quad 0 < t \leq 1$$

On cherche l'intersection avec le cercle, donc:

$$s^c + t(s^N - s^c) s^c = -\frac{g^T g}{g^T H g} g$$

$$\|p(t) - x^k\|^2 = \delta^2$$

$$(s^c + t(s^N - s^c))^T (s^c + t(s^N - s^c)) = \delta^2$$

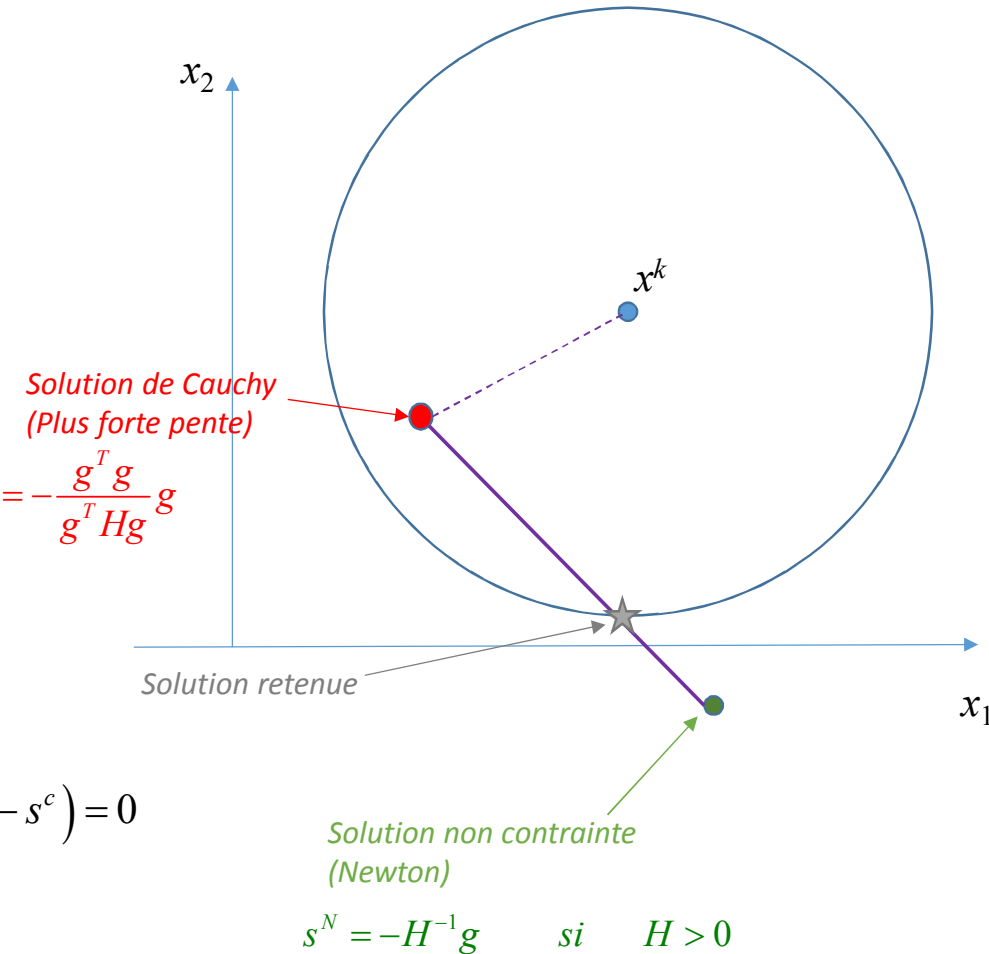
$$\left( (s^c)^T + t(s^N - s^c)^T \right) (s^c + t(s^N - s^c)) = \delta^2$$

$$(s^c)^T s^c - \delta^2 + t \left[ (s^N - s^c)^T s^c + (s^c)^T (s^N - s^c) \right] + t^2 (s^N - s^c)^T (s^N - s^c) = 0$$

$$a + bt + ct^2 = 0$$

2 racines, on choisit celle qui appartient au segment

$$s = s^c + t(s^N - s^c)$$



## Régions de confiance : Méthode « Dogleg »

### Algorithme:

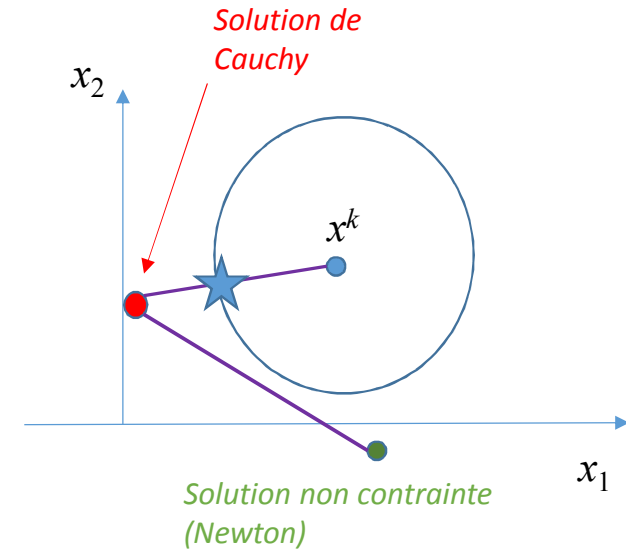
1) Vérifier que la méthode de Newton (non contrainte) améliore la solution

⇒ Permet de détecter les cas où le Hessien n'est pas défini positif

⇒ Dans ce cas, utiliser la méthode de Cauchy

Sol de Newton:  $s^N = -H^{-1} \cdot g$

Variation du modèle:  $q(x^{k+1}) - q(x^k) = g^T s^N + \frac{1}{2} g^T H g$



## Régions de confiance : Méthode « Dogleg »

### Algorithme:

2) Si  $q(x^{k+1}) - q(x^k) < 0$  alors

Si newton est dans le domaine (  $\|s^N\| \leq \delta$  )  $\Rightarrow$  on garde cette solution

Sinon

On calcule la solution de Cauchy (non contrainte) :  $s^c = -\frac{g^T g}{g^T H g} \cdot g$

Si Cauchy en dehors de la contrainte alors

on s'arrête en bordure de région

Sinon

$\Rightarrow$  Cauchy en dehors de la contrainte:

On choisit à l'intersection du cercle et du segment  $[s^c, s^N]$

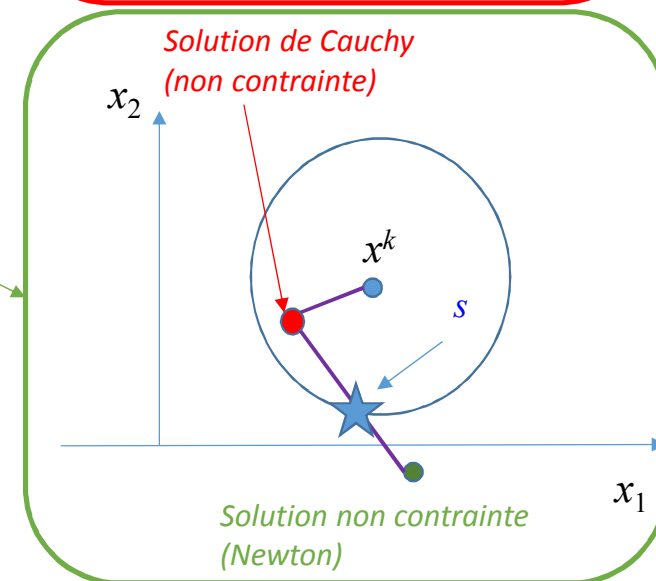
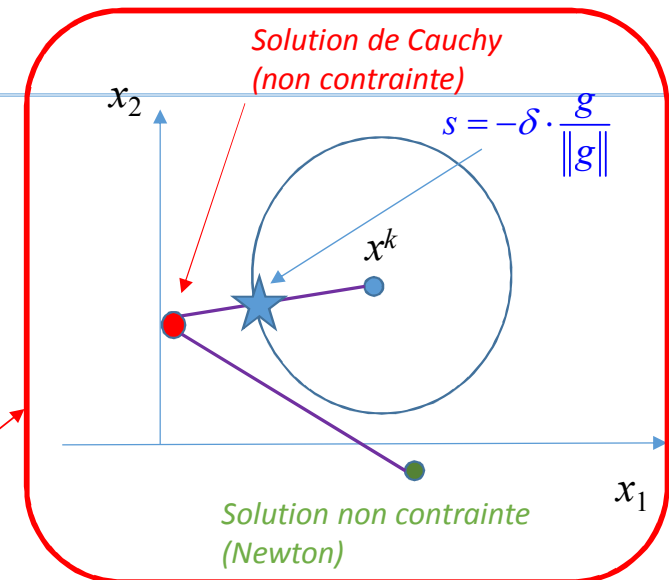
$$t = \frac{-b + \sqrt{b^2 - 4ac}}{2a}$$

$$s = s^c + t(s^N - s^c)$$

$$a = (s^N - s^c)^T (s^N - s^c)$$

$$b = (s^N - s^c)^T s^c + (s^c)^T (s^N - s^c)$$

$$c = (s^c)^T s^c - \delta^2$$



## Régions de confiance : Méthode « Dogleg »

```
function s=TR_Dodleg(g,H,Delta)
% Calcul de la solution uncontrainte avec Newton
sN=-H\g;

% Evaluation de q(xk+1)-q(x)
DiffQ=g'*sN+0.5*sN'*H*sN;

if DiffQ<0
    % La sol de Newton améliore la solution
    if norm(sN)<Delta
        % La solution est acceptée
        s=sN;
    else
        % La sol de Newton est en dehors de la région
        % On calcule la solution de Cauchy
        sC = -(g'*g)/(g'*H*g)*g;
        if norm(sC)>Delta
            % La sol de cauchy est en dehors
            % On limite la taille du pas dans la direction de descente
            s= -Delta*g/norm(g);
        else
            % La sol de Newton est dehors, celle de Cauchy dedans
            % On cherche l'intersection du cercle et du segment
            DiffS=sN-sC;
            a=DiffS'*DiffS;
            b=DiffS'*sC+sC'*DiffS;
            c=sC'*sC-Delta^2;
            t=(-b+sqrt(b^2-4*a*c))/2/a;
            s=sC+t*DiffS;
        end
    end
else
    % La solution de Newton n'améliore pas la solution: point de scelle ou
    % H>0
    % On utilise la méthode de Cauchy
    s=TR_Cauchy(g,H,Delta);
end
```

NB: l'algorithme principal reste inchangé

# Régions de confiance : Méthode « Dogleg »

## Algorithme

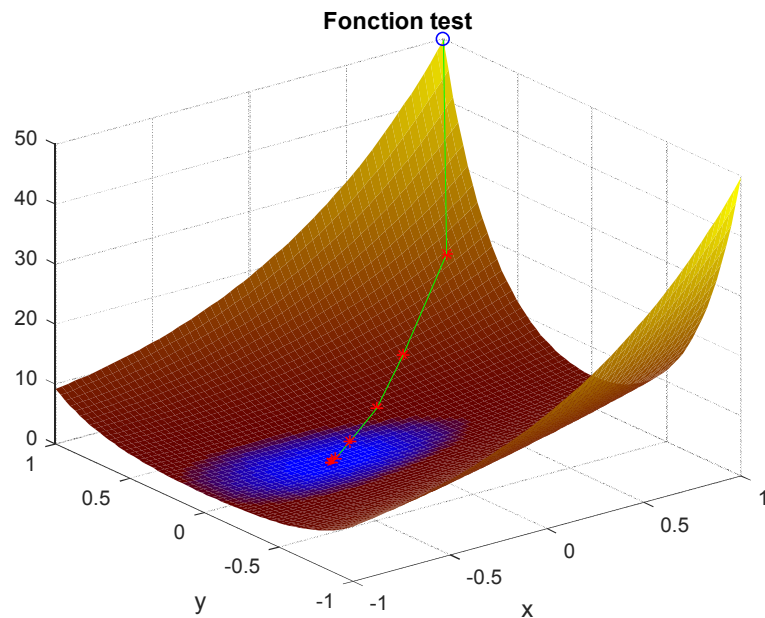
### Trust Region + Dogleg

9 itérations boucle principale  
25 appels de f

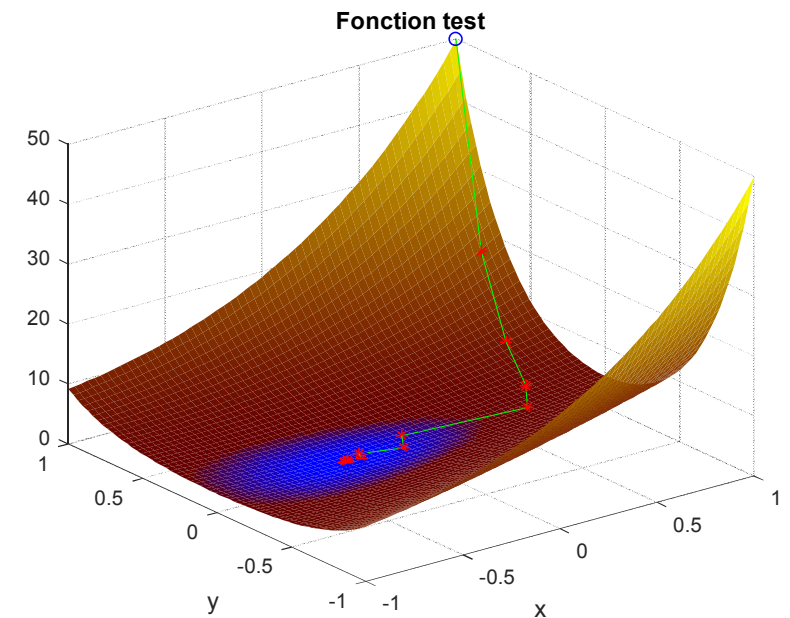
$$f(x, y) = e^{x-3y-\frac{1}{10}} + e^{x+3y-\frac{1}{10}} + e^{-x-\frac{1}{10}}$$

### Trust Region + Cauchy

17 itérations boucle principale  
49 appels de f



Dodleg: descent directement vers le minimum  
=> Avantage lié à la solution de Newton



Plus forte pente: descend vers la droite alors que le minimum est vers la gauche